

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6785738号  
(P6785738)

(45) 発行日 令和2年11月18日 (2020. 11. 18)

(24) 登録日 令和2年10月29日 (2020. 10. 29)

(51) Int. Cl. F I  
G O 6 F 12/00 (2006.01) G O 6 F 12/00 5 6 O F

請求項の数 9 (全 23 頁)

(21) 出願番号	特願2017-193447 (P2017-193447)	(73) 特許権者	390019839
(22) 出願日	平成29年10月3日 (2017. 10. 3)		三星電子株式会社
(65) 公開番号	特開2018-73452 (P2018-73452A)		S a m s u n g E l e c t r o n i c s
(43) 公開日	平成30年5月10日 (2018. 5. 10)		C o . , L t d .
審査請求日	令和2年9月9日 (2020. 9. 9)		大韓民国京畿道水原市靈通区三星路129
(31) 優先権主張番号	62/414, 426		129, S a m s u n g - r o , Y e o n
(32) 優先日	平成28年10月28日 (2016. 10. 28)		g t o n g - g u , S u w o n - s i , G
(33) 優先権主張国・地域又は機関	米国 (US)		y e o n g g i - d o , R e p u b l i c
(31) 優先権主張番号	15/425, 996	(74) 代理人	110000051
(32) 優先日	平成29年2月6日 (2017. 2. 6)		特許業務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	李 双 辰
早期審査対象出願			アメリカ合衆国, 93117, カリフ
			ォルニア州, ゴレタ, コレヒオ ロー
			ド, アパート #1302, 6510
			最終頁に続く

(54) 【発明の名称】 DRAM基盤のプロセッシングユニット

(57) 【特許請求の範囲】

【請求項 1】

D P U ( D y n a m i c R a n d o m A c c e s s M e m o r y P r o c e s s i n g U n i t ) であって、

少なくとも3以上である第1の所定の数のカラム及び3以上である第2の所定の数のローを有するアレイ内に配置された複数のDRAM基盤のコンピューティングセルを含む少なくとも1つのコンピューティングセルアレイを備え、

前記カラムの各々は、前記カラムの第1ロー及び第2ロー上で動作するロジック機能を提供し、前記カラムの第3ロー内に前記ロジック機能の結果を格納するように構成され、

前記少なくとも1つのコンピューティングセルアレイは、前記第1の所定の数の2倍である第3の所定の数のシフトラインを更に含み、

前記シフトラインの各々は、少なくとも1つの対応する第1のトランジスタを介して前記コンピューティングセルのカラムに連結され、

前記シフトライン及び前記対応する第1のトランジスタは、前記少なくとも1つのコンピューティングセルアレイ内において、選択されたカラムの2つのローのコンピューティングセルの内容を少なくとも2つのカラムで右方向又は左方向にシフトするように構成されることを特徴とするDPU。

【請求項 2】

前記第1の所定の数のカラム内に配置された少なくとも1つのDRAM基盤のメモリセルを含む少なくとも1つのデータセルアレイと、

10

20

前記コンピューティングセルのカラムの各々に連結されたセンスアンプと、を更に含み

、  
前記少なくとも1つのデータセルアレイのDRAM基盤のメモリセルのカラムの各々は、対応するコンピューティングセルアレイのカラムに対応し、

前記センスアンプは、前記コンピューティングセルのカラムの読出しビットラインに電氣的に連結される入力、及び前記コンピューティングセルのカラムの書込みビットラインに電氣的に連結される出力を含むことを特徴とする請求項1に記載のDPU。

【請求項3】

DPU(Dynamic Random Access Memory Processing Unit)であって、

少なくとも3以上である第1の所定の数のカラム及び3以上である第2の所定の数のローを有するアレイ内に配置された複数のDRAM基盤のコンピューティングセルを含む少なくとも1つのコンピューティングセルアレイと、

前記第1の所定の数のカラム及び少なくとも1つのロー内に配置された少なくとも1つのDRAM基盤のメモリセルを含む少なくとも1つのデータセルアレイと、を備え、

前記カラムの各々は、前記カラムの第1ロー及び第2ロー上で動作するロジック機能を提供し、前記カラムの第3ロー内に前記ロジック機能の結果を格納するように構成され、

前記少なくとも1つのコンピューティングセルアレイは、前記第1の所定の数の2倍である第3の所定の数のシフトラインを更に含み、

前記シフトラインの各々は、少なくとも1つの対応する第1のトランジスタを介して前記コンピューティングセルのカラムに連結され、

前記シフトライン及び前記対応する第1のトランジスタは、前記少なくとも1つのコンピューティングセルアレイ内において、選択されたカラムの2つのローのコンピューティングセルの内容を少なくとも2つのカラムで右方向又は左方向にシフトするように構成され、

前記少なくとも1つのデータセルアレイのDRAM基盤のメモリセルのカラムの各々は、対応するコンピューティングセルアレイのカラムに対応することを特徴とするDPU。

【請求項4】

DPU(Dynamic Random Access Memory Processing Unit)であって、

3以上である第1の所定の数のカラム及び少なくとも1つのロー内に配置された少なくとも1つのDRAM基盤のメモリセルを含む少なくとも1つのデータセルアレイと、

複数のDRAM基盤のコンピューティングセルを含む少なくとも1つのコンピューティングセルアレイと、

前記コンピューティングセルのカラムの各々に連結されたセンスアンプと、

前記コンピューティングセルの各々に電氣的に連結されるデコーダーと、を備え、

前記コンピューティングセルアレイの各々は、前記データセルアレイの各々に対応し、前記第1の所定の数のカラム及び3以上である第2の所定の数のローを有するアレイ内に配置され、

前記コンピューティングセルのカラムの各々は、前記コンピューティングセルのカラムの第1ロー及び第2ロー上で動作するロジック機能を提供し、前記コンピューティングセルのカラムの第3ロー内に前記ロジック機能の結果を格納するように構成され、

前記少なくとも1つのコンピューティングセルアレイは、前記第1の所定の数の2倍である第3の所定の数のシフトラインを更に含み、

前記シフトラインの各々は、少なくとも1つの対応する第1のトランジスタを介して前記コンピューティングセルのカラムに連結され、

前記シフトライン及び前記対応する第1のトランジスタは、前記少なくとも1つのコンピューティングセルアレイ内において、選択されたカラムの2つのローのコンピューティングセルの内容を少なくとも2つのカラムで右方向又は左方向にシフトするように構成され、

10

20

30

40

50

前記センスアンプの各々は、前記コンピューティングセルのカラム内の前記コンピューティングセルのカラムの読出しビットラインに電氣的に連結される入力、及び前記コンピューティングセルのカラムの前記コンピューティングセルのカラムの書込みビットラインに電氣的に連結される出力を含み、

前記デコーダーは、前記カラムのコンピューティングセルを選択するための命令に対応するD R A M基盤のアドレス信号を受信して前記カラムの第1ロー及び第2ロー上で前記ロジック機能を生じ、前記カラムの第3ロー内に前記ロジック機能の結果を格納することを特徴とするD P U。

【請求項5】

少なくとも1つのカラムの前記D R A M基盤のコンピューティングセルの各々は、3つのトランジスタ及び1つのキャパシターを含むD R A Mメモリセルを含むことを特徴とする請求項1、3、又は4に記載のD P U。

10

【請求項6】

少なくとも1つのカラムの前記D R A M基盤のコンピューティングセルは、N O Rロジック機能を提供することを特徴とする請求項5に記載のD P U。

【請求項7】

少なくとも1つのカラムの前記D R A M基盤のコンピューティングセルの各々は、1つのトランジスタ及び1つのキャパシターを含むD R A Mメモリセルを含むことを特徴とする請求項1、3、又は4に記載のD P U。

【請求項8】

20

前記D R A M基盤のコンピューティングセルの各々は、前記D R A M基盤のコンピューティングセルのビットラインに連結されたA L U ( A r i t h m e t i c L o g i c U n i t ) を更に含み、

前記A L Uは、前記ロジック機能を提供することを特徴とする請求項7に記載のD P U。

【請求項9】

前記A L Uは、N O Rロジック機能を提供することを特徴とする請求項8に記載のD P U。

【発明の詳細な説明】

30

【技術分野】

【0001】

本発明は、メモリシステムに関し、より詳細には、D R A M基盤のプロセッシングユニットに関する。

【背景技術】

【0002】

一般的に、G P U ( G r a p h i c s P r o c e s s i n g U n i t ) 及びT P U ( T e n s o r P r o c e s s i n g U n i t ) がディープラーニング ( d e e p l e a r n i n g ) プロセッシングのために使用される。ディープラーニングプロセッシングはG P U又はT P Uによって効率的に処理できない高度に並列化されたプロセッシングを含む。

40

【先行技術文献】

【特許文献】

【0003】

【特許文献1】米国特許第5,901,095号明細書

【特許文献2】米国特許第6,035,384号明細書

【特許文献3】米国特許第6,195,738号明細書

【特許文献4】米国特許第7,299,099号明細書

【特許文献5】米国特許第8,042,082号明細書

【特許文献6】米国特許第9,136,872号明細書

50

【特許文献7】米国特許第9,197,285号明細書  
 【特許文献8】米国特許第9,317,482号明細書  
 【特許文献9】米国特許第9,378,181号明細書  
 【特許文献10】米国特許出願公開第2007/0226567号明細書  
 【特許文献11】米国特許出願公開第2011/0302366号明細書  
 【特許文献12】米国特許出願公開第2012/0246380号明細書  
 【特許文献13】米国特許出願公開第2015/0089166号明細書  
 【特許文献14】米国特許出願公開第2015/0131383号明細書  
 【特許文献15】米国特許出願公開第2016/0147667号明細書  
 【特許文献16】米国特許出願公開第2016/0173102号明細書  
 【特許文献17】欧州特許出願公開第1193502号明細書  
 【特許文献18】欧州特許出願公開第2523352号明細書

10

【非特許文献】

【0004】

【非特許文献1】MATAM, Kiran et al., "Energy-Efficient Large-Scale Matrix Multiplication on FPGAs", 2013 International Conference on Reconfigurable Computing and FPGAs (ReConFig), December 9-11, 2013 (8 pages).

20

【非特許文献2】NODA, Hideyuki et al., "A cost-efficient high-performance dynamic TCAM with pipelined hierarchical searching and shift redundancy architecture", IEEE Journal of Solid-State Circuits, Vol. 40, No. 1, January 3, 2005 (10 pages).

【非特許文献3】SESHADRI, Vivek et al., "Fast Bulk Bitwise AND and OR in DRAM", IEEE Computer Architecture Letters, Vol. 14, No. 2, May 18, 2015 (6 pages).

【非特許文献4】SESHADRI, Vivek et al., "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization", MICRO-46, Proc. of 46th IEEE/ACM Int. Symp. on Microarchitecture, December 7, 2013, pp. 185-197.

30

【非特許文献5】WANG, Gesan et al., "TCAM-Based Forwarding Engine with Minimum Independent Prefix Set (MIPS) for Fast Updating", 2006 IEEE International Conference on Communications, June 11-15, 2006 (7 pages).

【発明の概要】

40

【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、DRAM基盤のプロセッシングユニット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DPU (Dynamic Random Access Memory Processing Unit) は、少なくとも1つのカラムを含むアレイ内に配置された複数のDRAM基盤のコンピューティングセルを含む少なくとも1つのコンピューティングセルアレイを備え、前記少なくとも1

50

つのカラムは、少なくとも3つのローのDRAM基盤のコンピューティングセルを含み、前記少なくとも3つのローのDRAM基盤のコンピューティングセルは、前記少なくとも3つのローの第1ロー及び第2ロー上で動作するロジック機能を提供し、前記少なくとも3つのローの第3ロー内に前記ロジック機能の結果を格納する。

#### 【0007】

上記目的を達成するためになされた本発明の他の態様によるDPU(Dynamic Random Access Memory Processing Unit)は、少なくとも1つのカラムを含むアレイ内に配置された複数のDRAM基盤のコンピューティングセルを含む少なくとも1つのコンピューティングセルアレイと、少なくとも1つのカラム内に配置された少なくとも1つのDRAM基盤のメモリセルを含む少なくとも1つのデータセルアレイと、を備え、前記少なくとも1つのカラムは、少なくとも3つのローのDRAM基盤のコンピューティングセルを含み、前記少なくとも3つのローのDRAM基盤のコンピューティングセルは、前記少なくとも3つのローの第1ロー及び第2ロー上で動作するロジック機能を提供し、前記少なくとも3つのローの第3ロー内に前記ロジック機能の結果を格納する。

#### 【0008】

上記目的を達成するためになされた本発明の更に他の態様によるDPU(Dynamic Random Access Memory Processing Unit)は、少なくとも1つのカラム内に配置された少なくとも1つのDRAM基盤のメモリセルを含む少なくとも1つのデータセルアレイと、少なくとも3つのローのDRAM基盤のコンピューティングセルを含む少なくとも1つのカラムを含むアレイ内に配置された複数のDRAM基盤のコンピューティングセルを含む少なくとも1つのコンピューティングセルアレイと、前記少なくとも3つのローのDRAM基盤のコンピューティングセルの読出しビットラインに電氣的に連結された入力、及び前記少なくとも3つのローのDRAM基盤のコンピューティングセルの書込みビットラインに電氣的に連結された出力を含むセンスアンプと、前記少なくとも3つのローのDRAM基盤のコンピューティングセルに電氣的に連結されたデコーダーと、を備え、前記少なくとも3つのローのDRAM基盤のコンピューティングセルは、前記少なくとも3つのローの第1ロー及び第2ロー上で動作するロジック機能を提供し、前記少なくとも3つのローの第3ロー内に前記ロジック機能の結果を格納し、前記デコーダーは、前記第1ロー及び第2ロー上で前記ロジック機能を生成し、前記第3ロー内に前記ロジック機能の結果を格納するために、前記少なくとも3つのローのDRAM基盤のコンピューティングセルを選択するための命令に対応するDRAM基盤のアドレス信号を受信する。

#### 【発明の効果】

#### 【0009】

本発明によれば、多様な動作に対してプログラム可能且つ再構成可能なDPUを提供することができる。また高性能、エネルギー効率、低価格のシステムを提供することができる。

#### 【図面の簡単な説明】

#### 【0010】

【図1】本発明の一実施形態によるDPUの一例を示すブロック図である。

【図2A】コンピューティングセルアレイ内のコンピューティングセルに対して使用される3つのトランジスタ及び1つのキャパシターによるDRAMコンピューティングセルトポグラフィの一例を示す図である。

【図2B】コンピューティングセルアレイ内のコンピューティングセルに対して使用される1つのトランジスタ及び1つのキャパシターによるDRAMコンピューティングセルトポグラフィの他の例を示す図である。

【図3】本発明の一実施形態によるイントラマットシフトアレイの一例を示す図である。

【図4A】本発明の一実施形態によるインターマットシフトアレイの一例を示す図である。

【図４Ｂ】本発明の一実施形態による左側のインターマットシフトに対する隣接するコンピューティングセルカラムで同様に配置された２つのコンピューティングセル間のインターマットシフトインターコネクション構成を概念的に示す図である。

【図４Ｃ】本発明の一実施形態による左側のインターマットシフトに対する隣接するコンピューティングセルカラムで同一でないように配置された２つのコンピューティングセル間のインターマットシフトインターコネクション構成を概念的に示す図である。

【図５】本発明の一実施形態によるインターマットフォワーディングアレイを示す図である。

【図６Ａ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

10

【図６Ｂ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

【図６Ｃ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

【図６Ｄ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

【図６Ｅ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

【図６Ｆ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

20

【図６Ｇ】本発明の一実施形態によるＤＰＵによって提供されるＮＯＲロジック基盤の動作を示す図である。

【図７】本発明の一実施形態による確率的データアレイを含むＤＰＵの一例を示すブロック図である。

【図８Ａ】多重化動作に変換される加算動作に対する確率的コンピューティング動作を示す図である。

【図８Ｂ】ＡＮＤロジック動作に変換される乗算動作に対する確率的コンピューティング動作を示す図である。

【図９】本発明の一実施形態によるＤＰＵを含むシステム構造を示す図である。

30

【発明を実施するための形態】

【００１１】

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。

【００１２】

本明細書で、多様な詳細な説明を本文の理解を提供するために提示する。しかし、このような詳細な説明無しに記載した本発明の思想が当業者によって容易に具現されることはよく理解される。他の例で、広く公知された方法、手続、構成、及び回路は本文を曖昧にしないために説明しない。

【００１３】

詳細な説明で“１つの実施形態”又は“一実施形態”を参照することは、実施形態に関連する特定の特徴、構造、又は特性が本文の少なくとも１つの実施形態に含まれることを意味する。即ち、本文の多様な箇所で使用する“１つの実施形態で”又は“一実施形態で”又は“１つの実施形態によって”又は類似な意味を有する他の表現は同一の実施形態を参照することを要求しない。更に、特定の特徴、構造、又は特性は適切な方式により１つ以上の実施形態で組合される。これに関連して、本明細書で使用するように、“例示的な”という単語は“例（example、instance、illustration）”として提供される”ということの意味する。本明細書で“例示的な”のように記述した実施形態は、他の実施形態に比べて必ずしも望ましいか又は有利であるものと考えてはならない。また、本文で言及する内容に従って、単数用語は複数の形態も含み、複数用語は単数形態も含む。構成図面を含む多様な図面は、説明の目的のためにのみ本文で言及し、正確な縮尺ではない。同様に、多様な波形及びタイミング図は説明の目的のためにのみ図示

40

50

する。例えば、一部の構成要素は明確性のために他の構成要素と比較して誇張して図示する。更に、適切に考えられる場合、参照符号は、対応する構成要素及び／又は類似の構成要素を示すために図面で反復する。

#### 【0014】

本明細書で使用する用語は、特定の実施形態のみを説明するものであり、本発明を制限しようとするものとして意図しない。本明細書で使用するように、文脈上で明確に異なっている意味しない限り、単数形態の“１つ”は複数の形態も含むものとして意図する。“構成される”、“構成されている”、“含む”、及び“含んでいる”の用語を本明細書で使用する場合、このような用語は、定まった特徴、整数、段階、動作、要素、及び／又は成分の存在を明示するが、１つ又はそれ以上の他の特徴、整数、段階、動作、要素、成分、及び／又はそれらのグループの追加又は存在を排除しない。“第１”、“第２”のような用語は、先に説明した構成に対するラベルとして使用され、別の定義が無い限り特定の順序（例えば、空間的、時間的、論理的、等）を意図するものではない。更に、同一の参照符号は、同一であるか若しくは類似の機能を有する部品、構成、ブロック、回路、ユニット、又はモジュールに関連する２つ以上の図面に亘って使用される。しかし、このような使用は、説明の簡易化のみのためであり、このような構成及びユニットの構成又は構造的な細部事項が全ての実施形態又は共通的に参照される部分／モジュールで同一なものとして意図せず、本発明の特定の実施形態のみを指称するための１つの手段である。

#### 【0015】

異なって定義しない限り、全ての用語（技術又は科学用語を含む）は本発明の装置と方法が属する分野で当業者に通常的に理解される同一な意味を有する。また、辞書的な意味として定義した用語は、関連する記述及び／又は本発明の説明の状況に従って解釈されなければならない、このように定義されない限り、理想的に理解されるか若しくは過度に形式的な意味として理解されてはならない。

#### 【0016】

本明細書に記述する発明は、多様な動作に対してプログラム可能（programmable）であり、再構成可能な（reconfigurable）DPU（DRAM（Dynamic Random Access Memory）based Processing Unit）を提供する。例えば、多様な動作は、加算、乗算、シフト、最大／最小（MIN/MAX）、及び比較（COMPARISON）等を含むが、本発明はこれに制限されない。一実施形態で、DPUは、３つのトランジスタ及び１つのキャパシター（３Ｔ１Ｃ）を含むDRAMプロセス及び構造に基づく。他の実施形態で、若干の変化と共に、DPUは、１つのトランジスタ及び１つのキャパシター（１Ｔ１Ｃ）を含むDRAMプロセス及び構造に基づく。従って、DPUは、特定のコンピューティングロジック回路（例えば、アダー（加算器）のような）を含まないが、高度の並列動作を使用するメモリセルを使用する計算を提供する。一実施形態で、DPUは、乗算（multiplication）動作がANDロジック動作に変換され、加算（addition）動作が多重化（multiplexing）動作に変換される確率的（stochastic）コンピューティングアレイを含む。

#### 【0017】

また、本明細書に記述する発明は、DPUをプログラムし、再構成するためのISA（Instruction Set Architecture）、コンパイラ、ドライバー、ライブラリ、フレームワーク拡張を有する環境（ecosystem）を含むシステム構造を提供する。

#### 【0018】

更に、本明細書に記述する発明は、データセンター及び／又はモバイルアプリケーションに適合するシステム構造を提供する。そして、システム構造は、GPU/ASIC（Application Specific Integrated Circuit）（TPU）/FPGA（Field-Programmable Gate Array）マシンラーニングアプリケーションに対して選択的に二進（binary）及び固定小数

点方式 (fixed point) の計算の両方に対するマシンラーニングアプリケーションのための PIM (Processor-in Memory) ソリューションを提供する。一実施形態で、本明細書に記述する発明は、高性能、エネルギー効率的、低価格のシステムを提供する。該当システムは、例えば二進加重ニューラルネットワーク (Binary Weight Neural Network) に対する加速化されたディープラーニングを提供する。

#### 【0019】

本明細書に記述する発明は、再構成及びプログラム可能であり、DRAM技術を利用して構成されるDPUに関連する。一実施形態で、DPUは、多様な動作(例えば、加算、乗算、整列、等)を遂行するように構成されるDRAM基盤のメモリセルアレイ及びDRAM基盤のコンピューティングセルアレイを含む。

10

#### 【0020】

DPUの内部構成は、サブアレイの複数のバンクに連結されたシステムバスを含む。一実施形態で、システムバスはサブアレイのHツリー連結バンクを提供するように構成される。各々のサブアレイはローカルコントローラを含み、各々の独立的なサブアレイは分離的に又は同時に活性化される。一実施形態で、DRAM基盤のセルは2つのアレイ(データセルアレイ及びコンピューティングセルアレイ)に区分される。一実施形態で、コンピューティングセルアレイはDRAM基盤のメモリセルで具現される。他の実施形態で、コンピューティングセルアレイはロジック回路を含むDRAM基盤のメモリセルで具現される。DPU内部構造は、またデータシフト及びデータ移動回路を含む。一実施形態で、確率的データ計算のために構成される第3DRAM基盤のセルアレイがある。

20

#### 【0021】

図1は、本発明の一実施形態によるDPU(DRAM(Dynamic Random Access Memory) based Processing Unit) 100の一例を示すブロック図である。DPU100は1つ以上のバンク(101a~101m)を含む。図1では、1つ以上のバンク(101a~101m)の中のバンク101a、101bのみを図示する。各バンク101は1つ以上のサブアレイ(102a~102n)を含む。図1では、1つ以上のサブアレイ(102a~102n)の中のサブアレイ102a、102bのみを図示する。また、各バンク101はバッファ103を含む。バッファ103は、個別サブアレイ102に連結され、システムバス104に連結される。バッファ103は、バンク102内の全体ロー(row)を読み出し、その後該当ローを再び同一のバンク又は他のバンクに書き込む。バッファ103は、また該当ローデータのコピーをサブアレイ102内の複数のマット(105a~105n)にブロードキャスト(broadcast)する。一実施形態で、バンク101及びシステムバス104はHツリー連結バンクを提供する。

30

#### 【0022】

各サブアレイ102は1つ以上のマット(又はレーン)105を含む。図1では、1つ以上のマット105の中のサブアレイ102aのマット105a~105nのみを図示する。各マット105は、データセルアレイ106、コンピューティングセルアレイ107、及びイントラマット(intra-mat)シフトアレイ108を含むDPU100の一領域である。マット105の例は点線109で囲まれるように図1で示される。各マット105は、データセルアレイデコーダー110、コンピューティングセルアレイデコーダー111、インターマット(inter-mat)シフトアレイ112、及びインターマットフォワードリングアレイ113を、隣接するマットと共有する。一実施形態で、データセルアレイデコーダー110、コンピューティングセルアレイデコーダー111、及びインターマットシフトアレイ112は、隣接するマット105の間にサブアレイコントローラ114と交互に物理的に配置される。一実施形態で、各デコーダー(110、111)は従来のDARMタイプのメモリデコーダーのように動作する。

40

#### 【0023】

一実施形態で、各マット105は通信的に(communicatively)サブア

50

レイコントローラ 114 に連結される。各サブアレイコントローラ 114 は他のサブアレイコントローラ 114 から独立して構成される。サブアレイコントローラ 114 はアドレス ( a d d r ) として命令を D R A M アドレスバスから受信する。アドレス ( 即ち、アドレス信号 ) に応答して、サブアレイコントローラ 114 は、データセルアレイ 106 及びコンピューティングセルアレイ 107 の中の少なくとも 1 つのアレイに出力としてデコーディングされたアドレスを提供する。即ち、サブアレイコントローラ 114 は、関連するデータセルアレイ 106 に対するデータセルアレイデコーダー 110 によってデコーディングされたソース / 目的地 ( d e s t i n a t i o n ) ( s r c / d s t ) アドレスを出力する。そして、コンピューティングセルアレイ 107 の場合、サブアレイコントローラ 114 は、コンピューティングセルアレイデコーダー 111 によってデコーディングされた動作 / 計算 ( o p / c a l c ) アドレスを出力する。また、サブアレイコントローラ 114 は、2 つ以上のサブアレイコントローラ 114 が組織化された方法で動作する D R A M バスからのアドレスとして命令を受信する。また、サブアレイコントローラ 114 はデータ移動回路を制御する。例えば、サブアレイコントローラ 114 は、イントラマツシフトアレイ 108、インターマツシフトアレイ 112、及びインターマツフォワードイングアレイ 113 を制御する。

#### 【 0 0 2 4 】

各データセルアレイ 106 は、少なくとも 1 つのカラム ( c o l u m n ) 及び少なくとも 1 つのロー ( r o w ) 内に配置される 1 つ以上の D R A M セルを含む。一実施形態で、データセルアレイ 106 は従来の D R A M セルアレイのように構成される。一実施形態で、データセルアレイ 106 は 2 K カラム及び 16 ローを含む。他の実施形態で、データセルアレイ 106 は 2 K カラムより少ないか又は多いカラムを含み、また 16 ローより少ないか又は多いローを含む。

#### 【 0 0 2 5 】

各コンピューティングセルアレイ 107 は、少なくとも 1 つのカラム及び少なくとも 1 つのロー内に配置される 1 つ以上のコンピューティングセルを含む。コンピューティングセルアレイ 107 内のカラム数はデータセルアレイ 106 内のカラム数と同一である。一実施形態で、コンピューティングセルアレイ 107 は 2 K カラム及び 16 ローを含む。他の実施形態で、コンピューティングセルアレイ 107 は 2 K カラムより少ないか又は多いカラムを含み、また 16 ローより少ないか又は多いローを含む。

#### 【 0 0 2 6 】

図 2 A は、コンピューティングセルアレイ 107 内のコンピューティングセルに対して使用される 3 つのトランジスタ及び 1 つのキャパシター ( 3 T 1 C ) による D R A M コンピューティングセルトポグラフィ ( t o p o g r a p h y ) 201 の一例を示す図である。図 2 A に示すように、ロー X 内の 3 T 1 C コンピューティングセルは第 1 トランジスタ T1 を含む。第 1 トランジスタ T1 は、書込みビットライン ( W r i t e B L ) に電氣的に連結されたソース端、キャパシター C1 の第 1 端及び第 2 トランジスタ T2 のゲート端の両側に電氣的に連結されたドレーン端、並びに書込みイネーブル ( W E N ) ラインに電氣的に連結されたゲート端を含む。キャパシター C1 の第 2 端はグラウンド ( g r o u n d ) ラインに電氣的に連結される。第 2 トランジスタ T2 は、グラウンドラインに電氣的に連結されたソース端、及び第 3 トランジスタ T3 のソース端に電氣的に連結されたドレーン端を含む。第 3 トランジスタ T3 は、ワードライン W L に電氣的に連結されたゲート端、及び読出しビットライン ( R e a d B L ) に電氣的に連結されたドレーン端を含む。3 T 1 C コンピューティングセルトポグラフィ 201 は、読出しビットライン ( R e a d B L ) に電氣的に連結された入力及び書込みビットライン ( W r i t e B L ) に電氣的に連結された出力を含むセンスアンプ ( S e n s e A m p l i f i e r : S A ) を含む。

#### 【 0 0 2 7 】

また、ロー Y 内のコンピューティングセル及びロー R 内のコンピューティングセルの両方は、ロー X 内のコンピューティングセルの配置と同様に、3 T 1 C D R A M 構成に配

10

20

30

40

50

置された3つのトランジスタ(T1~T3)及びキャパシターCを含む。図2Aに示した3つのコンピューティングセル及びセンスアンプ(SA)の一例は、NORロジック動作(即ち、'X NOR Y'ロジック動作)を提供するように構成される。該当NORロジック動作の結果はローRに格納される。3T1C DRAMコンピューティングセルの1つのカラムのみを図2Aに明示的に示したが、3T1Cコンピューティングセルが複数のカラム(例えば、2Kカラム)内に構成される等の他の実施形態が可能である。また、3つ以上のローが提供される他の実施形態が可能である。また、図2Aに示した3T1C DRAMコンピューティングセル構成はNORロジック動作を提供するが、3T1C DRAMコンピューティングセルトポグラフィ201のNORロジック動作は多様な機能的な動作を提供するために使用される。例えば、機能的な動作は、排他ノア(XNOR)、加算(ADD)、セレクト(SET)、MAX、SIGN、多重化(MUX)、CSA(Conditional Sum Addition)ロジック、乗算、ポップカウント(popcount)、COMPARE等を含む。但し、本発明はこれに制限されない。また、イントラマットシフトアレイ108及びインターマットシフトアレイ112はシフト機能を提供する。

10

#### 【0028】

図2Bは、図1のコンピューティングセルアレイ107内のコンピューティングセルに対して使用される1つのトランジスタ及び1つのキャパシター(1T1C)によるDRAMコンピューティングセルトポグラフィ(topography)202の他の例を示す図である。図2Bに示すように、1T1Cコンピューティングセルは第4トランジスタT4を含む。第4トランジスタT4は、キャパシターC2の第1端に電氣的に連結されたソース端、ビットラインBLに電氣的に連結されたドレーン端、及びワードラインWLに電氣的に連結されたゲート端を含む。キャパシターC2の第2端はグラウンドラインに電氣的に連結される。ビットラインBLはセンスアンプ(SA)の入力に電氣的に連結される。センスアンプ(SA)の出力は、多重化器(MUX)の第1入力、第5トランジスタT5のドレーン端、及びALU(Arithmetic Logic Unit)の入力に電氣的に連結される。多重化器(MUX)の出力はラッチ(LATCH)の入力に電氣的に連結される。第5トランジスタT5のソース端はラッチ(LATCH)の出力に電氣的に連結される。ALUの出力は多重化器(MUX)の第2入力に電氣的に連結される。図2Bで、第5トランジスタT5、多重化器(MUX)、ラッチ(LATCH)、及びALUは、それぞれコントローラ114から制御信号(CNTL1~NTL4)を受信する。一実施形態で、ALUはNOR機能を提供するように構成される。図2BのビットラインBLに電氣的に連結されたロジック回路はNORロジック動作を提供するが、ビットラインBLに電氣的に連結されたロジック回路(即ち、ALU)は他の機能動作(例えば、排他ノア(XNOR)、加算(ADD)、セレクト(SET)、MAX、SIGN、多重化(MUX)、CSA(Conditional Sum Addition)ロジック、乗算、ポップカウント(popcount)、COMPARE等)を提供する。但し、本発明はこれに制限されない。また、イントラマットシフトアレイ108及びインターマットシフトアレイ112はシフト機能を提供する。図2Bで1つの1T1Cコンピューティングセルのみを図示したが、複数のカラム及びローの1T1Cコンピューティングセルが提供される。

20

30

40

#### 【0029】

図2A及び図2Bから分かるように、DPUのコンピューティングセルは、特定の複雑なコンピューティングロジックを含まない。但し、代わりに、DPUのコンピューティングセルは、複数の多様なタイプの計算を遂行する機能(ability)を提供する再プログラム可能な性質(nature)を有する相対的に単純なトポグラフィを含む。また、DPUのボグラフィは、より多くの計算をより速くより効率的に遂行するためにメモリ構造に内在され、大量並列処理の長所を有するように配置される。

#### 【0030】

図3は、本発明の一実施形態によるイントラマット(intra-mat)シフト(s

50

h i f t) アレイ 108 の一例を示す図である。イントラマツトシフトアレイ 108 の記述 ( d e s c r i p t i o n ) を単純化するために、図 3 に示すように、4 つのカラムのコンピューティングメモリセルアレイ 107 の幅のマツト 105 を考慮する。イントラマツトシフトアレイ 108 は、アレイ内に配置された複数の第 6 トランジスタ T6 ( 図 3 では、1 つのトランジスタのみを T6 で表示する )、 $2^n$  シフトライン S L s ( n はマツト 105 内のコンピューティングセルのカラムである )、 $n+2$  シフトレフト ( l e f t ) コントロールライン S L c L s、2 シフトライト ( r i g h t ) コントロールライン S R c L s、及び n シフトマスクライン S M L s を含む。イントラマツトシフトアレイ 108 の第 6 トランジスタ T6 の一部は書込みビットライン ( W r i t e B L ) と  $2^n$  シフトライン S L s との間に電氣的に連結され、イントラマツトシフトアレイ 108 の他の第 6 トランジスタ T6 は読出しビットライン ( R e a d B L ) と  $2^n$  シフトライン S L s との間に連結される。このような第 6 トランジスタ T6 のゲートは  $n+2$  シフトレフトコントロールライン S L c L s 及び 2 シフトライトコントロールライン S R c L s に電氣的に連結される。イントラマツトシフトアレイ 108 の他の第 6 トランジスタ T6 は n シフトマスクライン S M L s と  $2^n$  シフトライン S L s との間に電氣的に連結される。イントラマツトシフトアレイ 108 のコントロールラインはマツト 105 に関連するサブアレイコントローラ 114 に電氣的に連結される。

#### 【0031】

コントロールライン ( S L c L s、S R c L s ) 上の適切な信号によって、イントラマツトシフトアレイ 108 は、マツト 105 内でデータをレフト ( l e f t ) シフト又はライト ( r i g h t ) シフトする。レフトシフトに対して、データは、サイン ( s i g n ) ビットで満たされ、1 つの動作毎に 1 ビット又は ( n - 1 ) ビット程シフトされる。n はマツト 105 当たりのカラム数である。ライトシフトに対して、命令による制御に従ってデータは 0 又は 1 で満たされる。或いは、データは  $2^0$ 、 $2^1$ 、...、 $2^{k-1}$ 、 $2^k$  マツト当たりのカラム数までシフトされる。 $2^k$  はカラム数である。

#### 【0032】

図 4 A は、本発明の一実施形態によるインターマツト ( i n t e r - m a t ) シフトアレイ 112 の一例を示す図である。インターマツトシフトアレイ 112 の記述 ( d e s c r i p t i o n ) を単純化するために、図 4 A ~ 図 4 C に示すように、マツト 105 が 2 つのカラムのコンピューティングメモリセルアレイ 107 の幅であるマツト 105 の構成を考慮する。即ち、各マツト 105 はコンピューティングメモリセル 107 a の第 1 カラム及びコンピューティングメモリセル 107 b の第 2 カラムを含む。インターマツトシフトアレイ 112 は、トランジスタ T112 a、T112 b、トランジスタ T112 c、T112 d、データシフトライン 112 e、112 f、及びインターマツトシフトコントロールライン I S L c L s を含む。マツト内で、トランジスタ T112 a は、コンピューティングメモリセル 107 a の第 1 カラムの読出しビットライン ( R e a d B L ) に電氣的に連結されたソース端、データシフトライン 112 e に電氣的に連結されたドレーン端を含む。トランジスタ T112 b は、コンピューティングメモリセル 107 b の第 2 カラムの読出しビットライン ( R e a d B L ) に電氣的に連結されたソース端、データシフトライン 112 f に電氣的に連結されたドレーン端を含む。データシフトライン ( 112 e、112 f ) はバッファ 103 に電氣的に連結される。バッファ 103 は図 4 A に図示していない。異なるマツトの間で、トランジスタ T112 c は、隣接マツト内のデータシフトライン 112 e にそれぞれ電氣的に連結されたソース端及びドレーン端を含む。トランジスタ T112 d は、隣接マツト内のデータシフトライン 112 f にそれぞれ電氣的に連結されたソース端及びドレーン端を含む。トランジスタ T112 c、T112 d のゲートはそれぞれ異なるインターマツトシフトコントロールライン I S L c L s のそれぞれに連結される。インターマツトシフトコントロールライン I S L c L s 上の適切な信号によって、インターマツトシフトアレイ 112 は、異なるマツトの間でデータをレフトシフト又はライトシフトする。インターマツトシフトアレイ 112 のコントロールラインはマツト 105 に関連するサブアレイコントローラ 114 に電氣的に連結される。

## 【0033】

図4Bは、本発明の一実施形態による左側のインターマツトシフトに対する隣接するコンピューティングセルカラム（マツト105a、105b）で同様に配置された2つのコンピューティングセル間のインターマツトシフトインターコネクション（interconnection）構成を概念的に示す図である。図4Bのインターコネクション構成を、利用されるインターコネクションノードによって概念的に図示し、利用されるインターコネクションノードを強調して図示する。例えば、トランジスタT112c、T112dは活性化され、これに従って導線経路が各々のトランジスタの間に形成される。従って、左側のコンピューティングセルカラム（マツト105a）と右側のコンピューティングセルカラム（マツト105b）との間でデータシフトライン（112e、112f）は連結される。トランジスタT112c、T112dのゲート端はアクティブなインターマツトシフトコントロールラインISLcLに電氣的に連結される。マツト105b内のトランジスタT112a、T112bは活性化され、従ってマツト105b内のコンピューティングセル107aの読出しビットライン（Read BL）は、マツト105bの左側であるマツト105a内のコンピューティングセル107aの書込みビットライン（Write BL）に電氣的に連結される。また、これに従って、マツト105b内のコンピューティングセル107bの読出しビットライン（Read BL）は、マツト105bの左側であるマツト105a内のコンピューティングセル107bの書込みビットライン（Write BL）に電氣的に連結される。

10

## 【0034】

20

図4Cは、本発明の一実施形態による左側のインターマツトシフトに対する隣接するコンピューティングセルカラム（105a、105b）で同一でないように配置された2つのコンピューティングセル間のインターマツトシフトインターコネクション（interconnection）構成を概念的に示す図である。図4Cのインターコネクション構成を、利用されるインターコネクションノードによって概念的に図示し、利用されるインターコネクションノードを強調して図示する。例えば、トランジスタT112c、T112dは活性化されて導線経路が各々のトランジスタの間に形成される。従って、左側のコンピューティングセルカラム（マツト105a）と右側のコンピューティングセルカラム（マツト105b）との間でデータシフトライン（112e、112f）は連結される。トランジスタT112c、T112dのゲート端はアクティブなインターマツトシフトコントロールラインISLcLに電氣的に連結される。マツト105aのトランジスタT112a、T112bは活性化され、従ってマツト105a内のコンピューティングセル107aの読出しビットライン（Read BL）は、マツト105aの右側であるマツト105b内のコンピューティングセル107aの書込みビットライン（Write BL）に電氣的に連結される。また、これに従って、マツト105a内のコンピューティングセル107bの読出しビットライン（Read BL）は、マツト105aの右側であるマツト105b内のコンピューティングセル107bの書込みビットライン（Write BL）に電氣的に連結される。

30

## 【0035】

図5は、本発明の一実施形態によるインターマツトフォワーディング（forwarding）アレイ113を示す図である。インターマツトフォワーディングアレイ113の説明を単純化するために、図5に示すように、2つのカラムのコンピューティングセルアレイ107の幅のマツト105の構成を考慮する。即ち、各マツト105はコンピューティングセル107aの第1カラム及びコンピューティングセル107bの第2カラムを含む。マツト105と共に、インターマツトフォワーディングアレイ113は、トランジスタT113a、T113b、トランジスタT113c、T113d、トランジスタT113e、T113f、 $2^n$ フォワーディングデータラインFDL、フォワーディングコントロールラインFCL、及び $2^m$ フォワーディングセクションラインFSLを含む。ここで、nはマツト内のコンピューティングセルのカラム数であり、mはセクション数である。トランジスタT113a、T113bのソース端はそれぞれコンピューティングセル107

40

50

aの第1カラムの書込みビットライン(W r i t e B L)及び読出しビットライン(R e a d B L)に電氣的に連結される。トランジスタT 1 1 3 a、T 1 1 3 bのドレーン端は第1データフォワーディングライン(F D L) 1 1 3 gに電氣的に連結される。トランジスタT 1 1 3 c、T 1 1 3 dのソース端はそれぞれコンピューティングセル1 0 7 bの第2カラムの書込みビットライン(W r i t e B L)及び読出しビットライン(R e a d B L)に電氣的に連結される。トランジスタT 1 1 3 c、T 1 1 3 dのドレーン端は第2データフォワーディングライン(F D L) 1 1 3 hに電氣的に連結される。トランジスタT 1 1 3 e、T 1 1 3 fのソース端はそれぞれトランジスタT 1 1 3 a、T 1 1 3 bのゲート端に電氣的に連結される。トランジスタT 1 1 3 e、T 1 1 3 fのドレーン端は両方とも同一のフォワーディングセクションラインF S Lに連結される。トランジスタT 1 1 3 e、T 1 1 3 fのゲート端はそれぞれ異なるフォワーディングコントロールラインF C L sに連結される。フォワーディングコントロールラインF C L s上の適切な信号によって、インターマットフォワーディングアレイ1 1 3は、マットの間でデータをフォワーディングする。インターマットフォワーディングアレイ1 1 3のコントロールラインは、相互間でデータがフォワーディングされるマット1 0 5に関連するサブアレイコントローラ1 1 4に電氣的に連結される。

#### 【0036】

図6A～図6Gは、本発明の一実施形態によるD P Uによって提供されるN O Rロジック基盤の動作を示す図である。図6A～図6Gで、第1オペランド(o p e r a n d)はローXに格納され、第2オペランドはローY又はローWに格納される。図6A～図6G内の矢印はコンピューティングセルの全体ローに対するN O Rロジック動作の入出力フローを示す。例えば、図6AのローXはローXのコンピューティングセルに格納されたオペランドの全体ローを示す。ローX内に格納されたオペランド及びローY内に格納されたオペランドのN O Rロジック動作の結果は結果ローR内に格納される。一実施形態で、ローX及びローYのオペランドは、例えば1 0 0個のカラム(例えば、 $x_1, x_2, \dots, x_{100}$ 及び $y_1, y_2, \dots, y_{100}$ )を含み、結果はローR(例えば、 $r_1, r_2, \dots, r_{100}$ )内に格納される。即ち、 $x_i \text{ NOR } y_i = r_i$ である。ここで、iはカラムインデックスである。他の実施形態で、ローXはロー内のコンピューティングセルの選択されたグループのみを示す。

#### 【0037】

図6Bはプレフィックス(p r e f i x) K o g g e - S t o n eアダー(加算器)に基づくNビット数に対するプールアダー動作を例示的に示す。図6Bで、第1NビットオペランドはローX内に格納され、第2NビットオペランドはローY内に格納される。図6B内に示した例示的な加算動作で、中間値(i n t e r m e d i a t e t e r m)( $G_0, P_0, G_1, P_1, G_2, P_2, \dots, G_{\log N + 1}, P_{\log N + 1}$ )が計算される。図6Bの最上のブロックはローX及びローYからの入力オペランドを利用して $G_0$ 及び $P_0$ を決定する5個の分離された動作を示す。第1動作で、最上のブロックはローXの逆(即ち、 $\sim X$ )を決定し、これはロー1に格納される。第2動作はローYの逆(即ち、 $\sim Y$ )を決定し、これはロー2に格納される。第3動作は「ローX NOR ローY」の動作を決定し、これはロー3に格納される。第4動作は「 $G_0 = \text{ロー1 NOR ロー2}$ 」を決定し、これはロー4に格納される。第5動作は「 $P_0 = \text{ロー3 NOR ロー4}$ 」を決定し、これはロー5に格納される。

#### 【0038】

図6Bの中間ブロックを参照すると、最上ブロックからの中間結果 $G_0, P_0$ が中間結果 $G_{i+1}, P_{i+1}$ を決定するために使用される。ここで、iはカラムインデックスである。即ち、図6Aの最上のブロックから決定された中間結果 $G_0, P_0$ が中間結果 $G_1, P_1$ を決定するために使用される。中間結果 $G_1, P_1$ は中間結果 $G_2, P_2$ を決定するために使用され、同様に中間結果 $G_{\log N + 1}, P_{\log N + 1}$ が決定される。図6Bの最下のブロックで、結果ローR 1及び結果ローR 2はそれぞれプールアダー(f u l l a d d e r)動作に対するキャリー(c a r r y)結果及び総合(s u m)結果を格

10

20

30

40

50

納する。

【0039】

図6Cは3T1C DRAMコンピューティングセルトポグラフィ201によって提供される例示的なセレクト動作を示す。ロー1はローXの逆（即ち、 $\sim X$ ）の中間結果を格納する。ロー2はローYの逆（即ち、 $\sim Y$ ）の中間結果を格納する。ロー3はローSの逆（即ち、 $\sim S$ ）の中間結果を格納する。ロー4は「ロー1 NOR ロー3」の中間結果を格納する。ロー5は「ロー2 NOR ローS」の中間結果を格納する。ロー6は「ロー4 NOR ロー5」の中間結果を格納する。ローRはロー6の逆の結果（ $S ? X : Y$ ）を格納する。

【0040】

図6Dは3T1C DRAMコンピューティングセルトポグラフィ201によって提供される他の例示的なセレクト動作を示す。ロー1はローXの逆（即ち、 $\sim X$ ）の中間結果を格納する。ロー2はローSの逆（即ち、 $\sim S$ ）の中間結果を格納する。ロー3は「ロー1 NOR ローS」の中間結果を格納する。ロー4はローXの逆（ $\sim X$ ）の中間結果を格納する。ローRは「ロー3 NOR ロー4」の結果（ $S ? X : \sim X$ ）を格納する。

【0041】

図6Eは3T1C DRAMコンピューティングセルトポグラフィ201によって提供される例示的なMAX/MIN動作を示す。ロー1はローYの逆（即ち、 $\sim Y$ ）の中間結果を格納する。ロー2はローX + ( $\sim Y + 1$ )の中間結果を格納する。ロー3は「 $C_{out} > n$ 」の中間結果を格納する。ロー4は「 $C_{out} ? X : Y$ 」の中間結果を格納する。ローRは「MAX(X : Y)」の結果を格納する。

【0042】

図6Fは3T1C DRAMコンピューティングセルトポグラフィ201によって提供される例示的な1ビット乗算動作を示す。ロー1は「ローX NOR ローW」の中間結果を格納する。ロー2は「ローX NOR ロー1」の中間結果を格納する。ロー3は「ローW NOR ロー1」の中間結果を格納する。結果ローRは「ロー2 NOR ロー3」の結果（即ち、「ローX XNOR ローW」の結果）を格納する。

【0043】

図6Gは3T1C DRAMコンピューティングセルトポグラフィ201によって提供される例示的な複数ビット乗算動作を示す。図6Gの上方のブロックで、ロー1はローWの逆（即ち、 $\sim W$ ）の中間結果を格納する。ロー2は $2^i$ 番目のレフトシフトされたローXの逆（即ち、 $\sim X < < 2^i$ ）の中間結果を格納する。ここで、 $i$ はインデックスである。ロー3は「ロー1 NOR ロー2」の中間結果（ $PP_i = \sim W \text{ NOR } \sim X < < 2^i$ ）を格納する。図6Gの下方のブロックで、ロー1は「ロー $PP_0$  SUM ロー $PP_i$ （ $PP_i$ ）」の中間結果を格納する。ロー2は「ロー2 NOR ロー $W_{sign}$ 」の中間結果を格納する。ローRは「 $X * W$ 」の結果を格納する。

【0044】

図7は、本発明の一実施形態による確率的（stochastic）データアレイ715を含むDPU700の一例を示すブロック図である。図1に示したDPU100の構成要素と同一な参照インジケータ（indicator）を有するDPU700の多様な構成要素は同様であり、このような同様の構成要素の記述はここで省略する。DPU700のサブアレイ102は、（実際の）データセルアレイ106、コンピューティングセルアレイ107、及びイントラマットシフトアレイ108と共に、確率的データアレイ715及びコンバーター確率（converter-to-stochastic）アレイ716を含む。

【0045】

確率的データアレイ715の各々は、少なくとも1つのカラム及び少なくとも1つのロー内に配置された1つ以上の確率的コンピューティングセルを含む。確率的データアレイ715内のカラム数はデータセルアレイ106及びコンピューティングセルアレイ107内のカラム数と同一である。一実施形態で、確率的データアレイ715は2Kカラム及び

10

20

30

40

50

16ローを含む。他の実施形態で、確率的データアレイ715は2Kカラムより少ないか又は多いカラム、或いは16ローより少ないか又は多いローを含む。確率的データアレイ715内で、'1'が存在する確率が使用され、2nビットはnビット値を示すために使用される。コンバーターツ確率アレイ716内の任意数生成器は実数(real number)を確率的な数に変換するために使用される。ポップカウント動作は確率的な数を再び実数に変換するために使用される。

#### 【0046】

確率的なコンピューティングアプローチを使用して、加算は多重化(multiplexing)動作に変換され、乗算はANDロジック動作に変換される。例えば、図8Aは、多重化動作に変換される加算動作に対する確率的コンピューティング動作を示す図であり、図8Bは、ANDロジック動作に変換される乗算動作に対する確率的コンピューティング動作を示す図である。確率的コンピューティングに対する従来の技術は巨大なメモリ容量を要求する。しかし、本明細書に記述した発明は高度に効率的な確率的コンピューティングを提供するために使用される。これはDRAM基盤のDPUが多くの並列AND動作及びマックス(MUX)動作を遂行するためである。本明細書に記述したDPUを使用する確率的コンピューティングは、またディープラーニングが一般的なアプリケーションである複雑な動作を加速化する。

#### 【0047】

図9は、本発明の一実施形態によるDPUを含むシステム構造900を示す図である。システム構造900は、ハードウェアレイヤー910、ライブラリ及びドライバレイヤー920、フレームワークレイヤー930、及びアプリケーションレイヤー940を含む。

#### 【0048】

ハードウェアレイヤー910は、本明細書に記述したDPUのような内装されたDPUを含むハードウェア装置及び/又は構成要素を含む。装置及び/又は構成要素の一実施形態は、1つ以上の内装されたDPUを含むPCIe装置911である。装置及び/又は構成要素の他の実施形態は、1つ以上の内装されたDPUを含むDIMM(Dual In-line Memory Module)912である。システム構造900のハードウェアレイヤー910はPCIe装置及び/又はDIMMに制限されず、ハードウェアレイヤー910はDPUを含むSOC(System On Chip)装置又は他のメモリタイプ装置を含むことは容易に理解される。ハードウェアレベル910で装置及び/又は構成要素内に内装されるDPUは、図1のDPU100及び/又は図7のDPU700と同様に構成される。他の実施形態で、DPUの特定コンピューティングセルアレイは、図2Aの3T1Cコンピューティングセルトポグラフィ201又は図2Bの1T1Cコンピューティングセルトポグラフィ202を含むように構成される。

#### 【0049】

システム構造900のライブラリ及びドライバレイヤー920は、DPUライブラリ921、DPUドライバ922、及びDPUコンパイラ923を含む。DPUライブラリ921は、アプリケーションレイヤー940で動作する多様なアプリケーションに対するハードウェアレイヤー910内のDPU内サブアレイの各々に対する最適のマッピング機能、リソース割当機能、及びスケジューリング機能を提供するように構成される。

#### 【0050】

一実施形態で、DPUライブラリ921は、移動、加算、乗算等の動作を含むフレームワークレイヤー930に対するハイレベルAPI(Application Programming Interface)を提供する。例えば、DPUライブラリ921は、また標準タイプのルーチン(routines)に対する具現を含む。標準タイプのルーチンに対する具現は、加速化されたディープラーニングプロセスに対して適用されるフォワード(forward)及びバックワード(backward)コンボリューション(convolution)、プーリング(pooling)、正規化(normalization)、及び活性化(activation)レイヤーを含む。但し、本発明はこ

10

20

30

40

50

れに制限されない。一実施形態で、DPUライブラリ921は、CNN(Convolution Neural Network)の全体コンボリューションレイヤーに対する計算をマッピングするAPI類似機能を含む。更に、DPUライブラリ921は、DPU上へのコンボリューションレイヤー計算のマッピングを最適化するためのAPI類似機能を含む。

#### 【0051】

また、DPUライブラリ921は、タスク(例えば、バッチ(batch)、出力チャンネル、ピクセル、入力チャンネル、コンボリューションカーネル)内の全ての個別又は複数の並列性(parallelism)を、チップ、バンク、サブアレイ、及び/又はマットレベルで該当DPU並列性にマッピングして、リソース割当を最適化するためのAPI類似機能を含む。更に、DPUライブラリ921は、性能(即ち、データ移動フロー)と電力消耗との間で均衡を維持(trade off)する初期化及び/又はランタイム時に最適のDPU構成を提供するAPI類似機能を含む。DPUライブラリ921によって提供される他のAPI類似機能はデザインノブ(knob)タイプ機能を含む。例えば、デザインノブタイプ機能は、バンク当たり活性化されたサブアレイの数、活性化されたサブアレイ当たりの入力機能マップの数、機能マップのパーティショニング(partitioning)、及び/又はコンボリューションカーネルの再使用スキームの設定を含む。また他のAPI類似機能は、各サブアレイに対して特定タスク(例えば、コンボリューションコンピューティング、チャンネル圧縮(sum up)、及び/又はデータディスパッチング(dispatching))を割当することによって、追加的なリソースの割当最適化を提供する。オペランドが整数と確率的数字との間で変換された場合、DPUライブラリ921は、精密度制限を満足させながらも、オーバーヘッドを最小化するAPI類似機能を含む。精密度が予想より低い場合、DPUライブラリ921は、確率的表現のための追加的なビットを使用して値を再び計算するか、又は他のハードウェア(例えば、CPU(Central Processing Unit))にタスクを分担(offload)するAPI類似機能を含む。

#### 【0052】

DPUライブラリ921は、またDPU内の活性化されたサブアレイをスケジュールすると同時にデータ移動をスケジュールして、データ移動がコンピューティング動作によって隠されるAPI類似機能を含む。

#### 【0053】

DPUライブラリ921の他の様相は追加的なDPU開発のための拡張インターフェイスを含む。一実施形態で、DPUライブラリ921は、NOR及びシフトロジックを利用して直接機能をプログラムして標準タイプ動作(例えば、加算、乗算、MAX/MIN等)及び他の動作が提供されるインターフェイスを提供する。拡張インターフェイスは、またインターフェイスを提供し、従ってDPUライブラリ921によって具体的に支援されない動作が、ライブラリ及びドライバレイヤー920で、SoCコントローラ(図示せず)、CPU/GPU構成要素、及び/又はCPU/TPU構成要素として分担される。DPUライブラリ921の他の様相は、DPUメモリがコンピューティングのために使用されない場合に、メモリの拡張としてDPUのメモリを使用するためのAPI類似機能を提供する。

#### 【0054】

DPUドライバ922は、DPUハードウェアレイヤーをシステムに集積するために、ハードウェアレイヤー910でのDPU、DPUライブラリ921、及びより高いレイヤーでのオペレーティングシステム(OS: Operating System)との間のインターフェイス連結を提供するように構成される。即ち、DPUドライバ922はDPUをシステムOS及びDPUライブラリ921に露出する。一実施形態で、DPUドライバ922は初期化時にDPUコントロールを提供する。一実施形態で、DPUドライバ922はDRAMタイプのアドレス又はDRAMタイプのアドレスのシーケンスの形態で命令をDPUに伝送し、DPUの内外へのデータ移動を制御する。DPUドライバ

ー 9 2 2 は、D P U - C P U 及び / 又は D P U - G P U 通信を処理すると共に多重 D P U 通信を提供する。

【 0 0 5 5 】

D P U コンパイラ 9 2 3 は、D P U ライブラリ 9 2 1 からの D P U コードを、D P U を制御するために D P U ドライバー 9 2 2 によって使用されるメモリアドレスの形態である D P U 命令にコンパイルする。D P U コンパイラ 9 2 3 によって生成された D P U 命令は、D P U 内の 1 つ及び / 又は 2 つのロー上で作動する単一命令（例えば、ベクトル命令、及び / 又は集合（g a t h e r e d）ベクトル、リードオン動作命令）である。

【 0 0 5 6 】

フレームワークレイヤー 9 3 0 は使いやすい（u s e r - f r i e n d l y）インターフェイスをライブラリ及びドライバレイヤー 9 2 0 並びにハードウェアレイヤー 9 1 0 に提供するように構成される。一実施形態で、フレームワークレイヤー 9 3 0 は、アプリケーションレイヤー 9 4 0 で広範囲なアプリケーションと互換可能な使いやすいインターフェイスを提供し、D P U ハードウェアレイヤー 9 1 0 をユーザーに透過的に（t r a n s p a r e n t）提供する。他の実施形態で、フレームワークレイヤー 9 3 0 は、定量化（q u a n t i t a t i o n）機能を、例えば T o r c h 7 タイプアプリケーションや T e n s o r F l o w タイプアプリケーションのような、既存の、従来の方法に追加するフレームワーク拡張を含む。但し、本発明はこれに制限されない。一実施形態で、フレームワークレイヤー 9 3 0 はトレーニングアルゴリズムに定量化機能を追加することを含む。他の実施形態で、フレームワークレイヤー 9 3 0 は、既存の割り算、乗算、平方根のバッチ正規化方法に対して、割り算、乗算、平方根の近似方法にシフトする置換を提供する。他の実施形態で、フレームワークレイヤー 9 3 0 は、ユーザーが計算のために使用するビット数を設定する拡張を提供する。他の実施形態で、フレームワークレイヤー 9 3 0 は、D P U ライブラリ及びドライバレイヤー 9 2 0 からフレームワークレイヤー 9 3 0 に多重 D P U A P I をラップ（w r a p）するための容量を提供する。従って、ユーザーは複数の G P U の使用と同様にハードウェアレイヤーで複数の D P U を使用することができる。フレームワークレイヤー 9 3 0 の他の機能は、ユーザーがハードウェアレイヤー 9 1 0 で D P U 又は G P U に機能を割り当てる。

【 0 0 5 7 】

アプリケーション 9 4 0 は、広範囲なアプリケーション（例えば、イメージタグ（t a g）プロセッシング、セルフドライビング／パイロットティング（p i l o t i n g）運送手段、アルファ碁タイプディープマインドアプリケーション、及び / 又は音声研究（s p e e c h r e s e a r c h）等）を含む。但し、本発明はこれに制限されない。

【 0 0 5 8 】

以上、本発明の実施形態について図面を参照しながら詳細に説明したが、本発明は、上述の実施形態に限定されるものではなく、本発明の技術的範囲から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。

【 符号の説明 】

【 0 0 5 9 】

1 0 0、7 0 0      D P U  
1 0 1 a ~ 1 0 1 m、1 0 5 a ~ 1 0 5 n      バンク  
1 0 2 a ~ 1 0 2 n      サブアレイ  
1 0 3      バッファ  
1 0 4      システムバス  
1 0 5、1 0 5 a ~ 1 0 5 n、1 0 9      マット（レーン）  
1 0 6      データセルアレイ  
1 0 7      コンピューティングセルアレイ  
1 0 7 a ~ 1 0 7 d      コンピューティングセル  
1 0 8      イントラマットシフトアレイ  
1 1 0      データセルアレイデコーダー

10

20

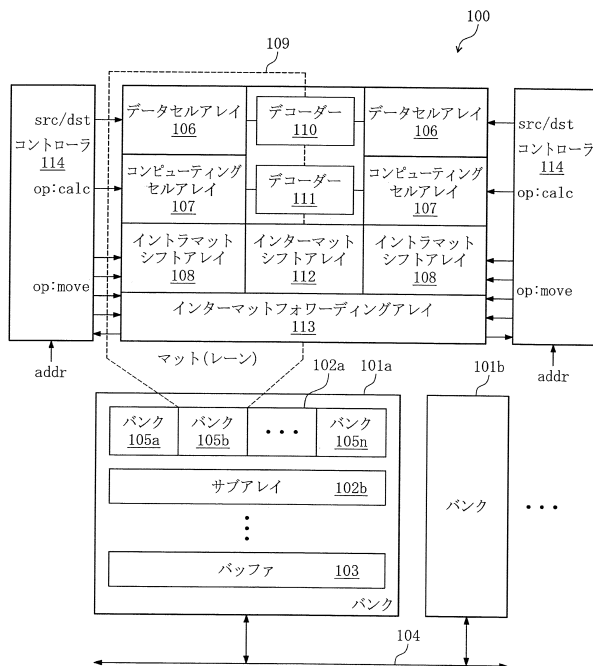
30

40

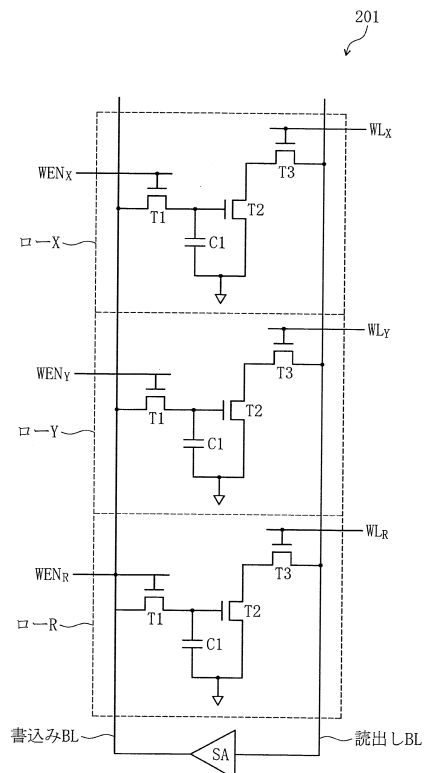
50

1 1 1	コンピューティングセルアレイデコーダー	
1 1 2	インターマツトシフトアレイ	
1 1 2 e、1 1 2 f	データシフトライン	
1 1 3	インターマツトフォワーディングアレイ	
1 1 3 g	第1データフォワーディングライン	
1 1 3 h	第2データフォワーディングライン	
1 1 4	サブアレイコントローラ	
2 0 1、2 0 2	D R A Mコンピューティングセルトポグラフィ	
7 1 5	確率的データアレイ	
7 1 6	コンバーターツ確率アレイ	10
9 0 0	システム構造	
9 1 0	ハードウェアレイヤー	
9 1 1	P C I e	
9 1 2	D I M M	
9 2 0	ライブラリ及びドライバレイヤー	
9 2 1	D P Uライブラリ	
9 2 2	D P Uドライバ	
9 2 3	D P Uコンパイラ	
9 3 0	フレームワークレイヤー	
9 4 0	アプリケーションレイヤー	20

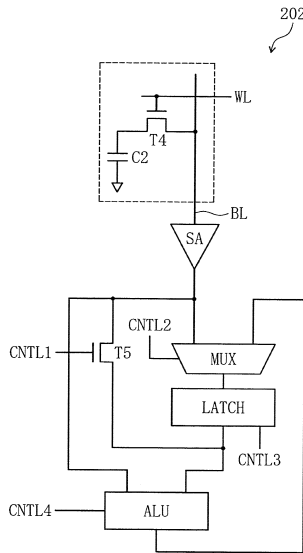
【図 1】



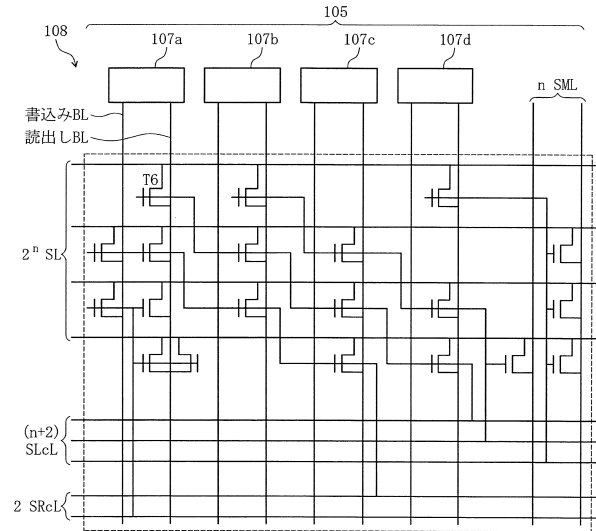
【図 2 A】



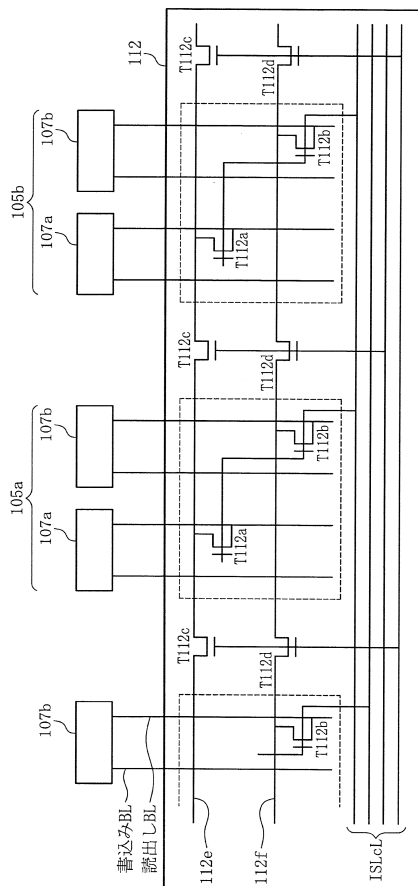
【図 2 B】



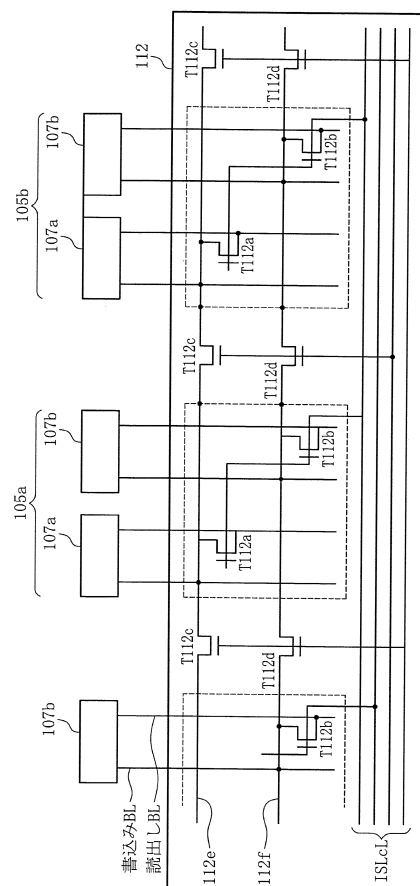
【図 3】



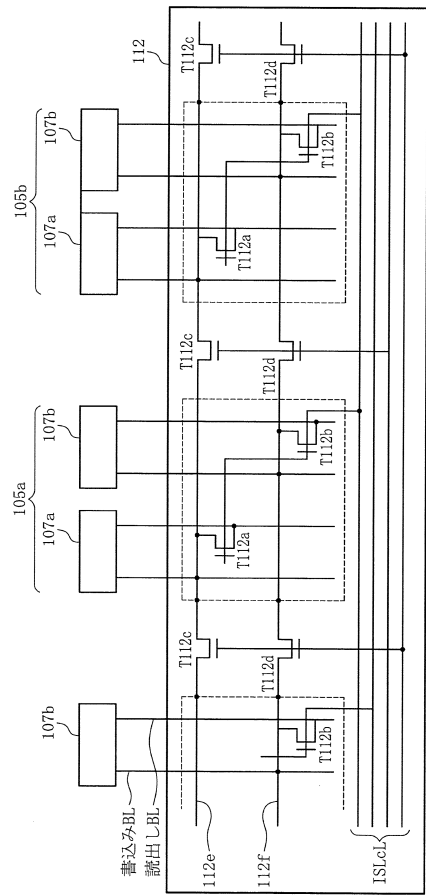
【図 4 A】



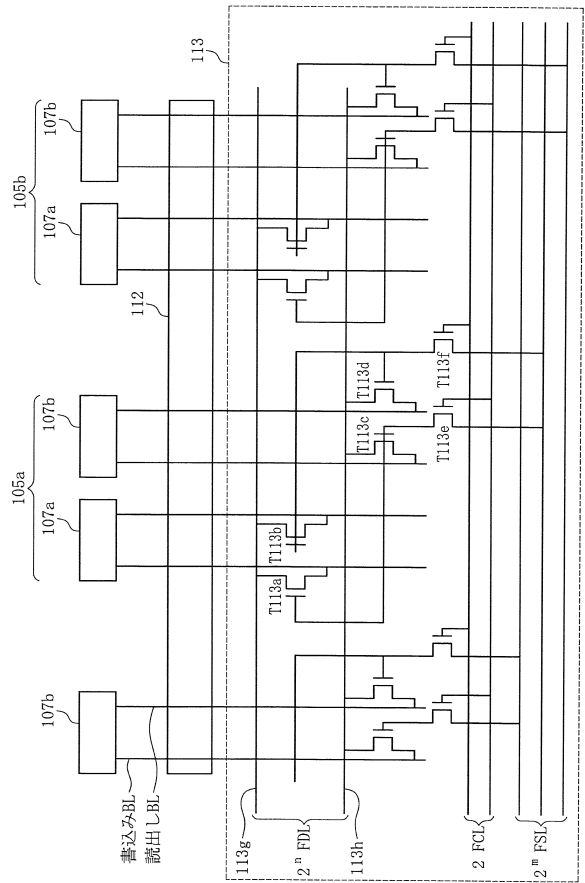
【図 4 B】



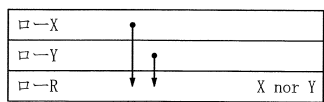
【図 4 C】



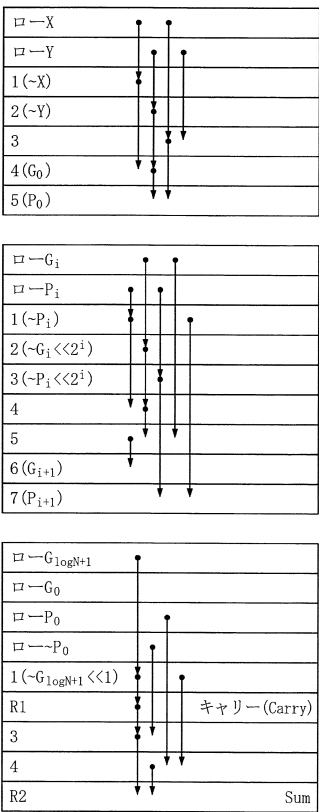
【図 5】



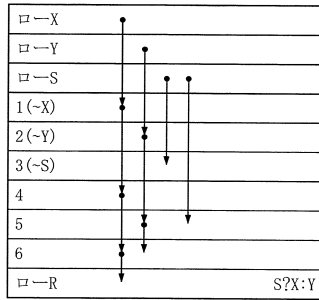
【図 6 A】



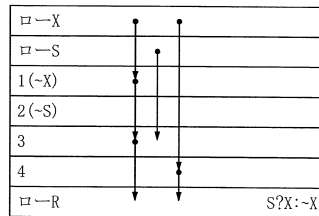
【図 6 B】



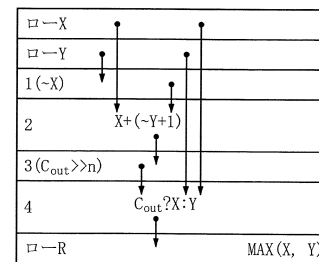
【図 6 C】



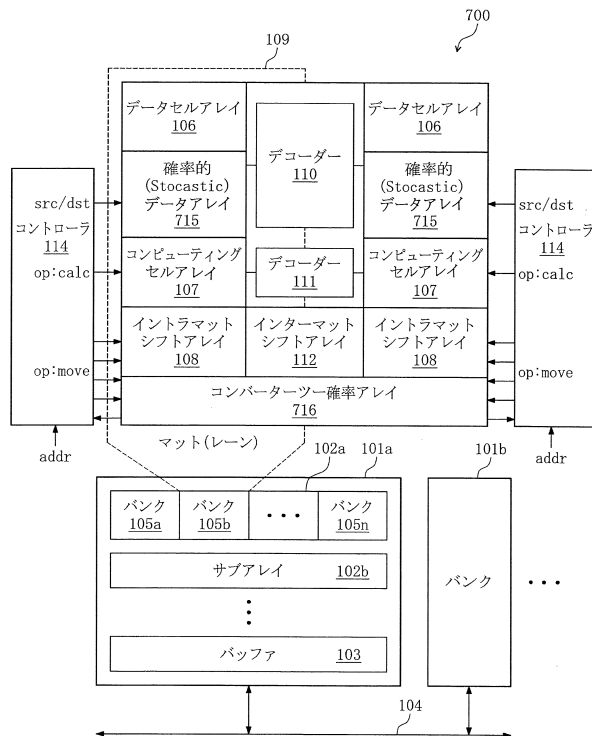
【図 6 D】



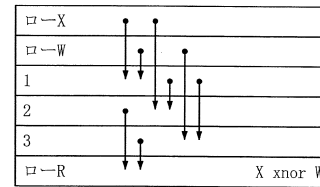
【図 6 E】



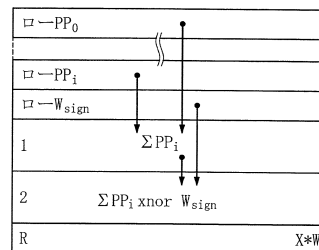
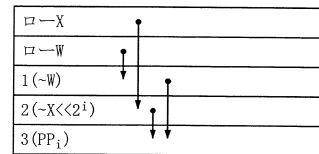
【図 7】



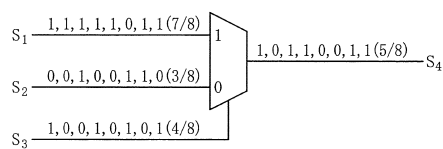
【図 6 F】



【図 6 G】

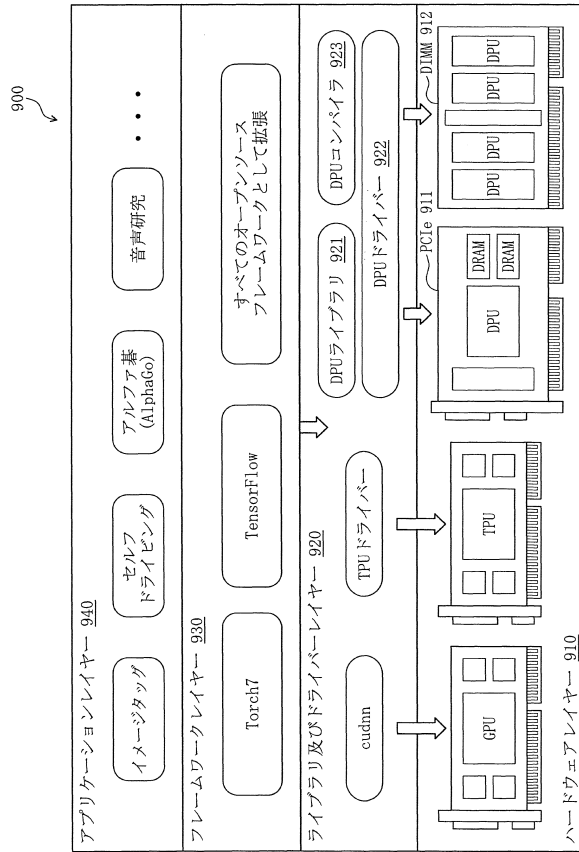


【図 8 A】



【図 8 B】





---

フロントページの続き

(72)発明者 牛 迪 民

アメリカ合衆国, 94087, カリフォルニア州, サニーベール, ホルトハウス テラス  
527

(72)発明者 マラディ, クリシュナ

アメリカ合衆国, 95135, カリフォルニア州, サン ノゼ, ロートレック ドライブ  
4196

(72)発明者 鄭 宏 忠

アメリカ合衆国, 95032, カリフォルニア州, ロス ガトス, カールトン アベニュー  
#6, 120

審査官 後藤 彰

(56)参考文献 特表2016-536733(JP,A)

特表2016-532919(JP,A)

米国特許出願公開第2009/0164789(US,A1)

米国特許第4068305(US,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00