

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2023年6月29日(29.06.2023)



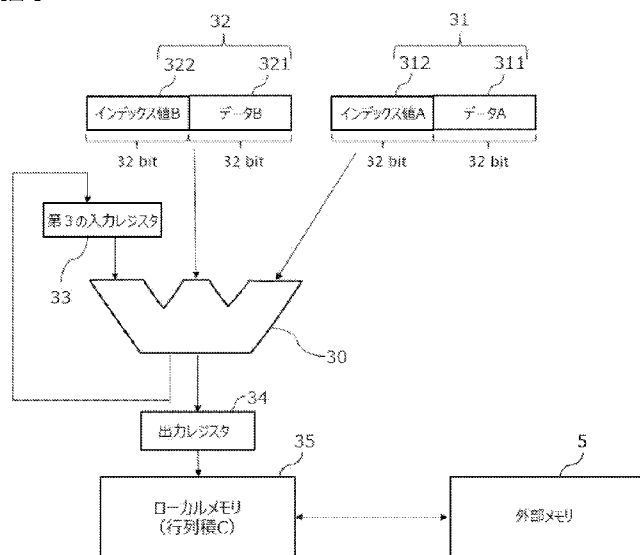
(10) 国際公開番号  
**WO 2023/120403 A1**

- (51) 国際特許分類:  
G06F 17/16 (2006.01) G06F 17/10 (2006.01)  
G06F 7/24 (2006.01)
- (21) 国際出願番号: PCT/JP2022/046353
- (22) 国際出願日: 2022年12月16日(16.12.2022)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:  
特願 2021-209979 2021年12月23日(23.12.2021) JP
- (71) 出願人: 国立大学法人奈良先端科学技術大学院大学(NATIONAL UNIVERSITY CORPORATION NARA INSTITUTE OF SCIENCE AND TECHNOLOGY) [JP/JP]; 〒6300192 奈良県生駒市高山町8916番地の5 Nara (JP).
- (72) 発明者: 中島 康彦(NAKASHIMA, Yasuhiko); 〒6300192 奈良県生駒市高山町8916番地の5 国立大学法人奈良先端科学技術大学院大学内 Nara (JP). 船井 遼太郎(FUNAI, Ryotaro); 〒6300192 奈良県生駒市高山町8916番地の5 国立大学法人奈良先端科学技術大学院大学内 Nara (JP).
- (74) 代理人: 弁理士法人グローバル知財(GLOBAL INTELLECTUAL PROPERTY); 〒6500021 兵庫県神戸市中央区三宮町3丁目7-6 神戸元町ユニオンビル9F Hyogo (JP).

(54) Title: CALCULATION UNIT INVOLVED IN MERGING AND SORTING AND PERFORMING SPARSE MATRIX COMPUTATION BY CGRA

(54) 発明の名称: CGRAによる疎行列計算とマージソートに関する演算ユニット

[図6]



- 5 External memory
- 33 Third input register
- 34 Output register
- 35 Local memory (matrix product C)
- 311 Data A
- 312 Index value A
- 321 Data B
- 322 Index value B

(57) Abstract: Provided is a unit that is for calculation by CGRA, that effectively uses a memory space, and that makes it possible to obtain a sparse matrix product by using less local memory. The calculation unit comprises a 3-input pipeline floating point sum-of-product calculator and a local memory. A sum-of-product calculator 30 is provided with: a first input register 31 and a second input register 32 that each are composed of a set of an index register and a data register; and a third input register 33 and an output register 34 that are data registers. In the first and second input registers of the



WO 2023/120403 A1

(81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CV, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IQ, IR, IS, IT, JM, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, WS, ZA, ZM, ZW.

(84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, CV, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, ME, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類 :

一 国際調査報告 (条約第21条(3))

sum-of-product calculator, index values A, B read from the local memory are stored in the respective index registers, data A, B read from the local memory are stored in the respective data registers, and a comparison is made between the index value A and the index value B. If the index values match with each other, a multiplication value of data A, B and a data value of the third input register are added, and the resultant value is stored in the output register of the sum-of-product calculator and also returns to the third input register.

(57) 要約 : メモリ空間を有効に利用し、少ないローカルメモリを用いて疎行列積を可能とするCGRAによる演算ユニットを提供する。演算ユニットは、3入力パイプライン浮動小数点積和演算器とローカルメモリから構成され、積和演算器30は、インデックスレジスタとデータレジスタを組にした第1の入力レジスタ31及び第2の入力レジスタ32、データレジスタである第3の入力レジスタ33及び出力レジスタ34を備える。積和演算器の第1及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値A、Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータA、Bが各データレジスタに格納され、インデックス値Aとインデックス値Bが比較され、値が一致する場合には、データA、Bの乗算値と第3の入力レジスタのデータ値が加算されて積和演算器の出力レジスタに格納されると共に、第3の入力レジスタに戻る。

## 明 細 書

発明の名称：

**C G R Aによる疎行列計算とマージソートに関する演算ユニット**

### 技術分野

[0001] 本発明は、2次元アレイ状に配置されたP E (Processing Element) を計算資源に持つC G R A (Coarse-grained reconfigurable architecture) における演算ユニットに関し、特に、疎行列計算とマージソートに関する演算ユニットに関するものである。

### 背景技術

[0002] 昨今、I o T (Internet of Things) など高い電力あたりの性能を持つデバイスが求められ、主にループ処理などをC P Uの代わりに実行するアクセラレータとして有用なC G R Aが注目されている。

C G R Aは、シストリックアレイを起源とするアーキテクチャーで、F P G A (Field Programmable Gate Array) と同様に回路を何度でも変更できる再構成可能なアーキテクチャーである。C G R Aは演算ユニット及びレジスタ等からなる基本ユニットとメモリから構成されることが多く、再構成の粒度がF P G Aではゲートレベルであるのに対し、C G R Aでは演算ユニットレベルである。そのため、C G R Aは、同等プロセスルールで同等機能を実現する場合には、F P G Aに比べ、動作周波数や回路面積で優位である。また、C G R Aは、パイプライン型垂直方向並列処理であり、演算ユニット群に毎サイクル供給するデータを抑えることが可能である。

[0003] ニューラルネットワークの計算などにおいて、大規模なパラメータ行列の計算をプロセッサで処理する際、計算量が多くなり処理負担が大きくなることから、処理負担を軽くすべく、主にF P G Aに関する疎行列ベクトル積を並列計算するアクセラレータにおいて、必要とするハードウェア資源を削減するアクセラレータが知られている（特許文献1を参照）。

特許文献1のアクセラレータは、データロード部が記憶部に記憶している

入力ベクトルの要素の各々に対して並列にアクセスする順序に従い、疎行列の非零要素の各々の配置にパディングを挿入する。しかしながら、特許文献1のアクセラレータにおける疎行列ベクトル積は、データバッファからの入力要素と重みパラメータ列の融合積和を行うものであり、これは密行列－疎行列の行列積となり、疎行列－疎行列の行列積に関するものではない。

[0004] また、データセンター向けに、CGRAを用いた密行列テンソル計算アクセラレータが実用化されているが、画像フィルタなど実際の畳み込みニューラルネットワーク(CNN)の計算を行うAI(Artificial Intelligence)アクセラレータでは、零要素の多い疎行列を取り扱うことが多い。しかしながら、大容量メモリを利用できるデータセンターでは、零要素を取り除くことなく実行しているのが実状である。

## 先行技術文献

### 特許文献

[0005] 特許文献1：特開2020-166368号公報

### 発明の概要

#### 発明が解決しようとする課題

[0006] メモリ容量に制約がある超小型のAIアクセラレータでは、零要素の多い疎行列を取り扱う場合、零要素を削除して、限りあるメモリ空間を有効に利用できる疎行列－疎行列の積を行うべきである。

かかる状況に鑑みて、本発明は、メモリ空間を有効に利用し、少ないローカルメモリを用いて疎行列－疎行列の積とマージソートを可能とするCGRAによる疎行列計算とマージソートに関する演算ユニットを提供することを目的とする。

#### 課題を解決するための手段

[0007] 上記課題を解決すべく、本発明の演算ユニットは、3入力パイプライン浮動小数点積和演算器とローカルメモリから構成される演算ユニットであって、積和演算器は、インデックスレジスタとデータレジスタを組にした第1の

入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備える。積和演算器の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納される。そして、インデックス値Aとインデックス値Bが比較され、値が一致する場合には、データAとデータBの乗算値と第3の入力レジスタのデータ値が加算されて積和演算器の出力レジスタに格納されると共に、第3の入力レジスタに戻る。

[0008] 上記構成によれば、CGRA型アクセラレータに、データパスを追加して、少ないローカルメモリを用いた疎行列-疎行列の行列積の演算を可能とする。具体的には、行列の要素番号であるインデックス値を上位32bit、行列の要素値であるデータを下位32bit格納する64bitデータ(8byte)をローカルメモリに格納し、疎行列-疎行列の行列積の演算を可能とする。また、2つの入力、疎行列であることに着目して、零要素を取り除き、データを圧縮して計算を行い、簡素な機構により、より高度な疎行列-疎行列の積を、疎行列-密行列を同じループで高効率に計算することを可能とする。CGRAはパイプライン型垂直方向であるため、垂直方向(縦方向)の圧縮はできないことから、水平方向(横方向)の零要素を取り除きデータを圧縮する。そして、データパスを追加することにより、少ないローカルメモリを用いた疎行列計算を途切れることなく可能とする。

[0009] 本発明の演算ユニットは、具体的には、疎行列Aと疎行列Bの行列積の演算において、インデックス値Aは疎行列の列番号で、インデックス値Bは疎行列の転置行列の列番号であり、各疎行列におけるゼロ値の要素は圧縮され、非ゼロ値の要素がインデックス値とデータを組としてローカルメモリに記憶される。

ここで、疎行列におけるゼロ値の要素が圧縮されるとは、ゼロ値のものを取り除き、非ゼロ値のものだけを残して、列番号を前に詰めることである。

インデックス値には、元の列番号が格納され、データには元の列番号の値が格納されている。

[0010] 本発明の演算ユニットにおいて、クロック信号のタイミング  $n+1$  で、積和演算器の第1段の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納される。

そして、クロック信号のタイミング  $n+2$  で、インデックス値Aとインデックス値Bが比較され、値が異なる場合には定数Cが積和演算器の第1段の出力データレジスタに格納され、値が一致する場合にはデータAとデータBの乗算が実行されて積和演算器の第1段の出力データレジスタに格納される。

クロック信号のタイミング  $n+3$  で、積和演算器の第2段の出力データレジスタに加算結果が伝搬され、クロック信号のタイミング  $n+4$  で、積和演算器の第3段の出力データレジスタに正規化済の積和演算結果が伝搬される。

そして、クロック信号のタイミング  $n+5$  で、積和演算器の第3段の出力データレジスタから積和演算器の第1段の第3の入力レジスタであるデータレジスタに戻される。

[0011] 本発明の演算ユニットは、入力ベースアドレスレジスタと入力オフセットアドレスレジスタを含む第1のアドレス計算機構および第2のアドレス計算機構を備え、第1および第2のアドレス計算機構は、各々の入力ベースアドレスレジスタおよび入力オフセットアドレスレジスタにアドレス情報をロードし、ローカルメモリから読み出したインデックス値Aとインデックス値Bを比較し、第1のアドレス計算機構は、 $A \leq B$ の場合に入力ベースアドレスレジスタに次要素参照に必要な値を加算し、第2のアドレス計算機構は、 $A \geq B$ の場合に入力ベースアドレスレジスタに次要素参照に必要な値を加算して、アドレス加算結果を入力ベースアドレスレジスタに戻す。

- [0012] 本発明の演算ユニットにおいて、クロック信号の第1のタイミングで、各々の入力ベースアドレスレジスタおよび入力オフセットアドレスレジスタにアドレス情報をロードし、クロック信号の第2のタイミングで、ローカルメモリから読み出したインデックス値Aとインデックス値Bを比較し、第1のアドレス計算機構は、 $A \leq B$ の場合に入力ベースアドレスレジスタに次要素参照に必要な値を加算し、第2のアドレス計算機構は、 $A \geq B$ の場合に入力ベースアドレスレジスタに次要素参照に必要な値を加算し、クロック信号の第3および第4のタイミングで、アドレス加算結果をデータ保持レジスタに伝搬させ、クロック信号の次のタイミングで、データ保持レジスタの内容を入力ベースアドレスレジスタに戻す。
- [0013] 本発明の演算ユニットにおいて、第1及び第2のアドレス計算機構は、ローカルメモリを複数空間に分割するアドレスマスク機構を備え、クロック信号の第1のタイミングで、入力オフセットアドレスレジスタに値をセットし、クロック信号の第2のタイミングで、入力オフセットアドレスレジスタの値をデータ保持レジスタに伝搬させ、クロック信号の第3のタイミングで、各々のアドレス加算結果と入力オフセットアドレスレジスタの値を加算し、クロック信号の第4のタイミングで、アドレスマスク機構により互いに異なる空間に分離してローカルメモリ参照のためのアドレスラッチに格納し、クロック信号の第5のタイミングで、ローカルメモリから読み出したインデックスとデータの組をメモリ出力レジスタに格納し、クロック信号の次のタイミングで、メモリ出力レジスタから読み出し、インデックス値として利用する。
- [0014] 本発明の第1の観点の演算方法は、3入力パイプライン浮動小数点積和演算器とローカルメモリから構成される演算ユニットを用いる演算方法であって、インデックスレジスタとデータレジスタを組にした第1の入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備える積和演算器を用い、下記1-1)～1-3)の各ステップを備える。

1-1) 積和演算器の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納されるステップ。

1-2) インデックス値Aとインデックス値Bが比較されるステップ。

1-3) インデックス値Aとインデックス値Bが一致する場合には、データAとデータBの乗算値と第3の入力レジスタのデータ値が加算されて積和演算器の出力レジスタに格納されると共に、第3の入力レジスタに戻すステップ。

[0015] 本発明の第1の観点の演算方法では、CGRA型アクセラレータに、データパスを追加して、少ないローカルメモリを用いた疎行列-疎行列の行列積の演算を行う。

行列の要素番号であるインデックス値を上位32bit、行列の要素値であるデータを下位32bit格納する64bitデータをローカルメモリに格納し、疎行列-疎行列の行列積の演算を行う。2つの入力、疎行列であることに着目して、零要素を取り除き、データを圧縮して計算を行い、疎行列-疎行列の積を、疎行列-密行列を同じスループットで高効率に計算する。CGRAはパイプライン型垂直方向であり、垂直方向(縦方向)の圧縮はできないことから、水平方向(横方向)の零要素を取り除きデータを圧縮し、データパスを追加することにより、少ないローカルメモリを用いた疎行列計算を途切れることなく行う。

[0016] 本発明の第1の観点の演算方法は、疎行列Aと疎行列Bの行列積の演算方法であって、下記1-4)、1-5)の各ステップを更に備えることが好ましい。

1-4) インデックス値Aは疎行列の列番号で、インデックス値Bは疎行列の転置行列の列番号で、各疎行列におけるゼロ値の要素が圧縮されるステップ。

1-5) 各疎行列における非ゼロ値の要素がインデックス値とデータを組と

してローカルメモリに記憶されるステップ。

[0017] 疎行列におけるゼロ値の要素の圧縮は、ゼロ値のものを取り除き、非ゼロ値のものだけを残して、列番号を前に詰め、インデックス値には元の列番号が格納され、データには元の列番号の値が格納される。

[0018] 本発明の第1の観点の演算方法は、下記1-6)~1-9)の各ステップを更に備えることが好ましい。

1-6) 各々の入力ベースアドレスレジスタおよび入力オフセットアドレスレジスタにアドレス情報をロードするステップ。

1-7) ローカルメモリから読み出したインデックス値Aとインデックス値Bを比較するステップ。

1-8) 第1のアドレス計算機構が、 $A \leq B$ の場合に入力ベースアドレスレジスタに次要素参照に必要な値を加算するステップ、或いは、第2のアドレス計算機構が、 $A \geq B$ の場合に入力ベースアドレスレジスタに次要素参照に必要な値を加算するステップ。

1-9) アドレス加算結果を入力ベースアドレスレジスタに戻すステップ。

[0019] 次に、本発明の第2の観点の演算方法は、3入力パイプライン浮動小数点積和演算器とローカルメモリから構成される演算ユニットを用いるマージソートの演算方法であって、インデックスレジスタとデータレジスタを組にした第1の入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備える積和演算器を用い、下記2-1)~2-5)の各ステップを備える。

2-1) 積和演算器の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納されるステップ。

2-2) インデックス値Aとインデックス値Bが比較されるステップ。

2-3) インデックス値Aとインデックス値Bのデータの読み出し結果の大小関係に従い、後続の演算ユニットにローカルメモリの2つの読み出しアド

レスと2つの読み出しデータを送り、アドレスとデータの各々の大小関係に従い、何れかのデータをストアするステップ。

2-4) マージソート全体のLog N段のうち、1段分のソート結果をローカルメモリにストアするステップ。

2-5) ローカルメモリのストア先アドレスを単調に増加するステップ。

[0020] 本発明の第2の観点の演算方法は、CGRA型アクセラレータに、データパスを追加して、少ないローカルメモリを用いて、データ全体を小さい又は大きい順にデータをソートしマージソートを行う。

ローカルメモリに格納された64bitデータの内、2つの入力の各インデックス値である各上位32bitのデータの読み出し結果の大小関係に従い、マージソートする。

### 発明の効果

[0021] 本発明のCGRAによる演算ユニットによれば、少ないローカルメモリを用いて疎行列-疎行列の積、マージソートが可能となるといった効果がある。

### 図面の簡単な説明

[0022] [図1]疎行列-疎行列の積をCGRAに実装する方法の説明図

[図2]疎行列の圧縮方法の説明図

[図3]従来の演算ユニットにおけるデータの足し合わせの説明図

[図4]本発明の演算ユニットにおけるデータ循環の説明図

[図5]本発明の演算ユニットの構成模式図

[図6]演算ユニットの構成図

[図7]アドレス計算機構の構成図

[図8]実施例1の演算ユニットの動作説明図(1)

[図9]実施例1の演算ユニットの動作説明図(2)

[図10]実施例1の演算ユニットの動作説明図(3)

[図11]実施例1の演算ユニットの動作説明図(4)

[図12]実施例1の演算ユニットの動作説明図(5)

[図13]実施例2の演算ユニットの動作説明図(1)

[図14]実施例2の演算ユニットの動作説明図(2)

[図15]実施例1と実施例2のアドレス計算部分の共通性についての説明図

[図16]疎行列積の演算フロー図

[図17]マージソートの演算フロー図

### 発明を実施するための形態

[0023] 以下、本発明の実施形態の一例を、図面を参照しながら詳細に説明していく。なお、本発明の範囲は、以下の実施例や図示例に限定されるものではなく、幾多の変更及び変形が可能である。

[0024] まず、疎行列-疎行列の積をCGRAに実装する方法について、図1を参照しながら説明する。ここでは、図1に示すように、4行8列の行列Aと8行4列の行列Bを掛け合わせて、4行4列の行列Cを計算する場合で説明する。行列Aと行列Bの行列積を計算する場合に、定義通りに計算すると、枠で囲むように、行方向と列方向を掛け合わせることになるが(図1(1)を参照)、これではメモリアクセスが順番にならないことから、行列Bの転置行列 $B^T$ を用いてメモリアクセスが順番になるようにする(図1(2)を参照)。

[0025] 行列A、行列Bが共に疎行列、すなわち、行列の要素の多くがゼロ値で構成される行列である場合には、非ゼロ値の要素のみを選択して計算することにより計算量を削減できる。例えば、図1(3)に示すように、行列Aの1行目の計算では、非ゼロ値のA00, A01だけを計算すればよく、その他のゼロ値の要素の計算はスキップできる。行列Aの1行目のA00, A01と掛け合わせる転置行列 $B^T$ の1~4行は、それぞれ、1行目はB00, B10、2行目はB01, B11, B21, B31、3行目はB42, B52, B62, B72、4行目はB63, B73となる。

[0026] そして、図1(4)に示す行列Aと転置行列 $B^T$ を圧縮する。

図2に示すように、圧縮は、行列Aと転置行列 $B^T$ における非ゼロ値の要素をインデックス値(行列の要素番号)とデータ(要素値)の組に置き換えて

、ゼロ値の要素を全て取り除く。元の行列のサイズを維持するため、残りのインデックス値とデータの組は、インデックス値として負の値（例えば-1）をセットする。

[0027] 図3に示す従来の演算では、例えば、行列Aの要素A00が、行列Bの要素B00, B01, B02, B03と掛け合わされたデータが、縦方向に流れ、後段の行列Aの要素A01が、行列Bの要素B10, B11, B12, B13と掛け合わされたデータと足し合わされ、同様に、行列Aの要素A02, A03が、それぞれ、行列Bの要素B20~B23, B30~B33と掛け合わされたデータと足し合わされ、最後に、行列Cの要素C00, C01, C02, C04となる。このように、データの足し合わせは、縦方向と同じように演算ユニットを並べるが、図1(4)に示す行列Aと転置行列B<sup>T</sup>を圧縮したものでは、縦方向と同じように演算ユニットをつなぎ合わせて実行することは不可能である。

[0028] そのため、図4に示すとおり、行列Aと転置行列B<sup>T</sup>を圧縮したものは、1つの演算ユニットの内部でインデックス値を参照して疎行列のデータを循環させて、データの足し合わせを行う必要がある。

[0029] 以下では、疎行列Aと疎行列Bの転置行列B<sup>T</sup>とを圧縮したもの（圧縮済み疎行列）を外部メモリから取り出して、1つの演算ユニットの内部でインデックス値を参照して疎行列のデータを循環させて、データの足し合わせを行う演算ユニットについて説明を行う。

[0030] 図5は、本発明の演算ユニットの構成模式図を示す。本発明の演算ユニットは、第1のアドレス計算機構1と第2のアドレス計算機構と演算ユニット3から構成され、それぞれが内部に持つローカルメモリ（15, 25, 35）は外部メモリ5と接続されている。

第1のアドレス計算機構1の演算器10では、入力ベースアドレスAと、疎行列Aと転置行列B<sup>T</sup>のインデックス値が入力され、各々のインデックス値を比較し、 $B \geq A$ である場合には、インデックス値Aのアドレスが増加され、疎行列Aの次要素の値が参照される。次要素のインデックス値は、インデ

ックス値とデータが組となり、例えば、64 bit (= 32 bit + 32 bit) 離れて記憶されているとすると、+8 byte アドレスを増加する。一方、各々のインデックス値を比較し、 $B \geq A$  でない場合には、インデックス値 A のアドレスはそのまま維持し、アドレスは増加しない。

また、演算器（加算器）11では、演算器10が出力するアドレスに、入力オフセットアドレスAの値が加算される。ローカルメモリ15には、疎行列Aのインデックス値とデータの組が記憶されている。メモリ出力レジスタ16は、ローカルメモリ15から読み出したインデックス値とデータの組を保持する。

[0031] 第2のアドレス計算機構2の演算器20では、入力ベースアドレスBと、疎行列Aと転置行列 $B^T$ のインデックス値が入力され、各々のインデックス値を比較し、 $B \leq A$  である場合には、インデックス値Bのアドレスが増加され、転置行列 $B^T$ の次要素の値が参照される。次要素のインデックス値は、インデックス値とデータが組となり、例えば、64 bit (= 32 bit + 32 bit) 離れて記憶されているとすると、+8 byte アドレスを増加する。一方、各々のインデックス値を比較し、 $B \leq A$  でない場合には、インデックス値Bのアドレスはそのまま維持し、アドレスは増加しない。

また、演算器（加算器）21では、演算器20が出力するアドレスに、入力オフセットアドレスBの値が加算される。ローカルメモリ25には、転置行列 $B^T$ のインデックス値とデータの組が記憶されている。メモリ出力レジスタ26は、ローカルメモリ25から読み出したインデックス値とデータの組を保持する。

[0032] 演算ユニット3は、3入力パイプライン浮動小数点積和演算器30とローカルメモリ35から構成される。演算器30は、第1及び第2のアドレス計算機構から出力された第1の入力レジスタ31と第2の入力レジスタ32と前回のプロセスの出力結果（第3の入力レジスタ33）とから3入力を行う。ローカルメモリ35は、外部メモリ5と接続される。

[0033] 図5に示す本発明の演算ユニットでは、下記（1）～（5）の4段パイプ

ライン処理の5つのデータフローが連携して途切れなく演算を継続する。4段パイプライン処理は、4段階で一つの命令の処理が終わるようにし、それぞれの処理が独立に、かつ、同時に処理され、並行して動作できる処理である。

[0034] (1) 浮動小数点積和演算器とローカルメモリのデータフロー（4段パイプライン処理）

クロック信号のタイミング  $n+1$  で、積和演算器の第1段の第1の入力レジスタ及び第2の入力レジスタに、ローカルメモリから読み出したインデックス値  $A$  とインデックス値  $B$  が各インデックスレジスタに格納され、ローカルメモリから読み出したデータ  $A$  とデータ  $B$  が各データレジスタに格納される。

クロック信号のタイミング  $n+2$  で、インデックス値  $A$  とインデックス値  $B$  が比較され、値が異なる場合には定数  $C$  が積和演算器の第1段の出力データレジスタに格納され、値が一致する場合にはデータ  $A$  とデータ  $B$  の乗算が実行されて積和演算器の第1段の出力データレジスタに格納される。

クロック信号のタイミング  $n+3$  で、積和演算器の第2段の出力データレジスタに加算結果が伝搬される。クロック信号のタイミング  $n+4$  で、積和演算器の第3段の出力データレジスタに正規化済の積和演算結果が伝搬される。先頭へ戻るデータ累算リングが構成される。

[0035] (2) 第1のアドレス計算機構のアドレス計算のデータフロー（4段パイプライン処理）

クロック信号の第1のタイミングで、入力ベースアドレスレジスタ  $A$  および入力オフセットアドレスレジスタ  $A$  にアドレス情報をロードする。

クロック信号の第2のタイミングで、ローカルメモリから読み出したインデックス値  $A$  とインデックス値  $B$  を比較し、入力ベースアドレスレジスタ  $A$  に次要素参照に必要な値（例えば、 $+8\text{ byte}$ ）を加算するか、或いは、そのままの値を維持する（ $+0\text{ byte}$ ）。

クロック信号の第3および第4のタイミングは通過し、先頭へ戻るアドレ

ス累算リングが構成される。

[0036] (3) 第2のアドレス計算機構のアドレス計算のデータフロー（4段パイプライン処理）

クロック信号の第1のタイミングで、入力ベースアドレスレジスタBおよび入力オフセットアドレスレジスタBにアドレス情報をロードする。

クロック信号の第2のタイミングで、ローカルメモリから読み出したインデックス値Aとインデックス値Bを比較し、入力ベースアドレスレジスタBに次要素参照に必要な値（例えば、+8 byte）を加算するか、或いは、そのままの値を維持する（+0 byte）。

クロック信号の第3および第4のタイミングは通過し、先頭へ戻るアドレス累算リングが構成される。

[0037] (4) 第1のアドレス計算機構のデータ参照フロー（4段パイプライン処理）

クロック信号の第2のタイミングで、ローカルメモリから読み出したインデックス値A、Bを比較し、入力ベースアドレスレジスタAに+0 byte又は+8 byteを加算出力し、

クロック信号の第3のタイミングで、入力オフセットアドレスレジスタAを加算し、

クロック信号の第4のタイミングで、アドレスマスク機構により互いに異なる空間に分離してローカルメモリ参照のためのアドレスラッチに格納し、

クロック信号の第5のタイミングで、ローカルメモリから読み出したインデックスとデータの組をメモリ出力レジスタに格納し、先頭に戻る。

[0038] (5) 第2のアドレス計算機構のデータ参照フロー（4段パイプライン処理）

クロック信号の第2のタイミングで、ローカルメモリから読み出したインデックス値A、Bを比較し、入力ベースアドレスレジスタBに+0 byte又は+8 byteを加算出力し、

クロック信号の第3のタイミングで、入力オフセットアドレスレジスタB

を加算し、

クロック信号の第4のタイミングで、アドレスマスク機構により互いに異なる空間に分離してローカルメモリ参照のためのアドレスラッチに格納し、

クロック信号の第5のタイミングで、ローカルメモリから読み出したインデックスとデータの組をメモリ出力レジスタに格納し、先頭に戻る。

[0039] (演算ユニットの動作)

図6を参照して、本発明の演算ユニットの動作について説明する。

本発明の演算ユニットは、3入力パイプライン浮動小数点積和演算器30とローカルメモリ35から構成され、積和演算器30は、インデックスレジスタ312とデータレジスタ311を組にした第1の入力レジスタ31、インデックスレジスタ322とデータレジスタ321を組にした第2の入力レジスタ32、データレジスタである第3の入力レジスタ33、及び、出力レジスタ34を備える。64bitデータで構成される第1の入力レジスタ31には、上位32bitのインデックスレジスタ312に疎行列Aの要素番号が、下位32bitのデータレジスタ311に要素値が格納される。また、64bitデータで構成される第2の入力レジスタ32には、上位32bitのインデックスレジスタ322に疎行列Bの転置行列B<sup>T</sup>の要素番号が、下位32bitのデータレジスタ321に要素値が格納される。

[0040] 具体的には、積和演算器30の第1の入力レジスタ31及び第2の入力レジスタ32には、図示しないローカルメモリ(行列積A, 行列積B)から読み出したインデックス値Aとインデックス値Bが各インデックスレジスタ(312, 322)に格納され、データAとデータBが各データレジスタ(311, 321)に格納される。

そして、インデックス値Aとインデックス値Bが比較され、値が一致する場合には、データAとデータBの乗算値と第3の入力レジスタ33のデータ値(前回のプロセスのデータ結果)が加算されて積和演算器30の出力レジスタ34に格納され、かつ、第3の入力レジスタ33に戻る。インデックス値Aとインデックス値Bの値が一致しない場合には、データAとデータBの

乗算をしないでスキップし、第3の入力レジスタ33のデータ値（前回のプロセスのデータ結果）をそのまま出力レジスタ34に出力する。

[0041] このように、1つの演算ユニットの内部で2つのインデックス値A、Bを参照し、疎行列のデータを出力から入力に戻し循環させて、疎行列積のデータの足し合わせを行っている。

そして、出力レジスタ34の内容をローカルメモリ35の行列積Cの各要素のアドレス領域に格納する。ローカルメモリ35と外部メモリ5は接続され、行列積Cの各要素は外部メモリ5に保存される。

[0042] （アドレス計算機構の動作）

次に、図7を参照して、本発明のアドレス計算機構の動作について説明する。

本発明のアドレス計算機構は、入力ベースアドレスレジスタ（18，28）と入力オフセットアドレスレジスタ（19，29）を含む第1のアドレス計算機構および第2のアドレス計算機構を備える。

[0043] 第1のアドレス計算機構は、クロック信号の第1のタイミングで、入力ベースアドレスレジスタA18および入力オフセットアドレスレジスタA19にアドレス情報をロードする。クロック信号の第2のタイミングで、ローカルメモリ15から読み出したインデックス値Aとインデックス値Bを比較し、 $A \leq B$ の場合に入力ベースアドレスレジスタA18に次要素参照に必要な値（+8 byte）を加算する。クロック信号の第3および第4のタイミングで、アドレス加算結果をデータ保持レジスタA14に伝搬させ、クロック信号の次のタイミングで、データ保持レジスタA14の内容を入力ベースアドレスレジスタA18に戻す。

[0044] 第2のアドレス計算機構は、クロック信号の第1のタイミングで、入力ベースアドレスレジスタB28および入力オフセットアドレスレジスタB29にアドレス情報をロードする。クロック信号の第2のタイミングで、ローカルメモリ25から読み出したインデックス値Aとインデックス値Bを比較し、 $A \geq B$ の場合に入力ベースアドレスレジスタB28に次要素参照に必要な

値 (+ 8 b y t e) を加算する。クロック信号の第3および第4のタイミングで、アドレス加算結果をデータ保持レジスタ B 2 4 に伝搬させ、クロック信号の次のタイミングで、データ保持レジスタ B 2 4 の内容を入力ベースアドレスレジスタ B 2 8 に戻す。

[0045] 第1及び第2のアドレス計算機構は、それぞれ、ローカルメモリを複数空間に分割するアドレスマスク機構 (1 0 1, 2 0 1) を備える。

クロック信号の第1のタイミングで、入力オフセットアドレスレジスタ (1 9, 2 9) に値をセットする。クロック信号の第2のタイミングで、入力オフセットアドレスレジスタ (1 9, 2 9) の値をデータ保持レジスタ (図示せず) に伝搬させる。クロック信号の第3のタイミングで、各々のアドレス加算結果と入力オフセットアドレスレジスタの値を加算する。

[0046] クロック信号の第4のタイミングで、アドレスマスク機構により互いに異なる空間に分離してローカルメモリ参照のためのアドレスラッチに格納する。クロック信号の第5のタイミングで、ローカルメモリから読み出したインデックス値とデータの組をメモリ出力レジスタ (1 6, 2 6) に格納する。クロック信号の次のタイミングで、メモリ出力レジスタ (1 6, 2 6) から読み出し、インデックス値として利用する。

### 実施例 1

[0047] 4クロック単位で動作するCGRAによる疎行列計算に関する演算ユニットの一実施形態について、具体的に説明する。以下の説明で参照する図8～12において、図中に数字が付されているが、下記説明におけるレジスタ番号と対応する (第nのレジスタは、図中のマルで囲んだ数字nに対応する)。

[0048] 演算ユニットと複数のアドレス生成器 (アドレス計算機構) とローカルメモリの組から構成される基本ユニットを複数連結する。各アドレス生成器の第1ステージは、ベースアドレスを保持する第1のレジスタ (図7における入力ベースアドレスレジスタ 1 8, 2 8) の値と、自身が最終ステージにおいて生成するアドレス計算結果をアドレス生成器毎に保持する第9のレジス

タ（図7におけるデータ保持レジスタ14，24）の値と、ローカルメモリの2カ所から読み出された2つのデータを各々保持する第3および第4のレジスタ（図7におけるメモリ出力レジスタ16，26）の値とを入力する。

[0049] 初回のループイタレーションでは、各アドレス生成器の第1ステージは第1のレジスタ（入力ベースレジスタ18，28）の値を加算の第1オペランド、定数0を加算の第2オペランドとして加算結果を第2のレジスタに格納する。

2回目以降のループイタレーションでは、各アドレス生成器の第1ステージはアドレス生成器毎の第9のレジスタ（データ保持レジスタ14，24）の値を第1オペランド、第3および第4のレジスタ（メモリ出力レジスタ16，26）の各上位bitの比較結果によりセットされる定数を第2オペランドとして加算結果を各アドレス生成器の第2のレジスタに格納する。

[0050] 各アドレス生成器の第2ステージは、第2のレジスタの値を加算の第1オペランド、オフセットアドレスを保持する第5のレジスタ（入力オフセットアドレスレジスタ19，29）の値を加算の第2オペランドとして加算を行い、加算結果を第6のレジスタに格納し、同時に、第2のレジスタの内容をそのまま第7のレジスタに格納する。

[0051] 各アドレス生成器の第3ステージは、第6のレジスタの値をローカルメモリ参照用マスク操作の第1オペランド、マスク値を第2オペランドとしてマスク操作を行い、マスク結果を各アドレス生成器の第8のレジスタに格納し、同時に、第7のレジスタの内容をそのまま各アドレス生成器の第9のレジスタ（データ保持レジスタ14，24）に格納する。

[0052] 各アドレス生成器が生成し第8のレジスタに格納するアドレス情報を用いて、ローカルメモリの複数個所を参照し、ロード結果を第3および第4のレジスタ（メモリ出力レジスタ16，26）に各々格納すると同時に、各々演算ユニットの境界に配置されるレジスタファイル中の第10および第11のレジスタ（図5，図6における第1及び第2の入力レジスタ31，32）に格納する。

[0053] 第10および第11のレジスタ（第1及び第2の入力レジスタ31, 32）の各下位32bitのデータA, Bを演算オペランドとして演算ユニット30に入力する。

## 実施例 2

[0054] 4クロック単位で動作するCGRAによるマージソートに関する演算ユニットの一実施形態について、具体的に説明する。以下の説明で参照する図13~15において、図中に数字が付されているが、下記説明におけるレジスタ番号と対応する（第nのレジスタは、図中のマルで囲んだ数字nに対応する）。

マージソートとは、データ全体を小さい又は大きい順にデータをソートしマージする演算である。図13に示すとおり、実施例1の疎行列計算と同様に計算を行い、第10および第11のレジスタ（図6における第1及び第2の入力レジスタ31, 32）の各上位32bitのデータの読み出し結果の大小関係に従い、1次元配列内の異なるアドレスA, Bの片側を更新してマージソートしていく。なお、実施例1で説明した疎行列計算に関する演算ユニットにおいて、上位32bitのインデックス値を比較するのと同様に、本実施例のマージソートに関する演算ユニットでも上位32bitをデータとして比較する。そして、マージソートの場合、第10および第11のレジスタの各下位32bitは、データの付随情報として、単にストアされるだけである。

[0055] 本実施例のマージソートに関する演算ユニットについても、実施例1の疎行列計算と同様に、演算ユニットと複数のアドレス生成器（アドレス計算機構）とローカルメモリの組から構成される基本ユニットを複数連結する。各アドレス生成器の第1ステージは、ベースアドレスを保持する第1のレジスタ（図7における入力ベースアドレスレジスタ18, 28）の値と、自身が最終ステージにおいて生成するアドレス計算結果をアドレス生成器毎に保持する第9のレジスタ（図7におけるデータ保持レジスタ14, 24）の値と、ローカルメモリの2カ所から読み出された2つのデータを各々保持する第

3および第4のレジスタ（図7におけるメモリ出力レジスタ16，26）の値とを入力する。

そして、初回のループイタレーションでは、各アドレス生成器の第1ステージは第1のレジスタ（入力ベースレジスタ18，28）の値を加算の第1オペランド、定数0を加算の第2オペランドとして加算結果を第2のレジスタに格納する。なお、2回目以降のループイタレーションでは、各アドレス生成器の第1ステージはアドレス生成器毎の第9のレジスタ（データ保持レジスタ14，24）の値を第1オペランド、第3および第4のレジスタ（メモリ出力レジスタ16，17）の各上位bitの比較結果によりセットされる定数を第2オペランドとして加算結果を各アドレス生成器の第2のレジスタに格納する。

[0056] 各アドレス生成器の第2ステージは、第2のレジスタの値を加算の第1オペランド、オフセットアドレスを保持する第5のレジスタ（入力オフセットアドレスレジスタ19，29）の値を加算の第2オペランドとして加算を行い、加算結果を第6のレジスタに格納し、同時に、第2のレジスタの内容をそのまま第7のレジスタに格納する。

次に、各アドレス生成器の第3ステージは、第6のレジスタの値をローカルメモリ参照用マスク操作の第1オペランド、マスク値を第2オペランドとしてマスク操作を行い、マスク結果を各アドレス生成器の第8のレジスタに格納し、同時に、第7のレジスタの内容をそのまま各アドレス生成器の第9のレジスタ（データ保持レジスタ14，24）に格納する。

そして、各アドレス生成器が生成し第8のレジスタに格納するアドレス情報を用いて、ローカルメモリの複数個所を参照し、ロード結果を第3および第4のレジスタ（メモリ出力レジスタ16，26）に各々格納すると同時に、各々演算ユニットの境界に配置されるレジスタファイル中の第10および第11のレジスタ（図5，図6における第1及び第2の入力レジスタ31，32）に格納する。

[0057] 第10および第11のレジスタ（第1及び第2の入力レジスタ31，32

) の各上位 32 bit のデータの読み出し結果の大小関係に従い、後続の演算ユニット 30 にアドレス A, B と、2 つの読み出しデータを送り、アドレスとデータの各々の大小関係に従い、何れかのデータをストアすることにより、ソート全体の Log N 段のうち、1 段分のソート結果をローカルメモリにストアする。ストア先アドレスは単調増加していく。

同時に後続の演算ユニット 30 が、前回の実行結果を前段のローカルメモリから読み出して、Log N 段の次段以降を担当することから、全体として、ローカルメモリをダブルバッファとするパイプライン実行が可能となるのである。

このように、本演算ユニットは、実施例 1 に示した疎行列積の演算の後、後続の演算ユニットにアドレス A, B と、2 つの読み出しデータを送り、アドレスとデータの各々の大小関係に従い、何れかのデータをストアさせることにより、マージソートの演算として使用できることがわかる。

[0058] 図 14 は、4 クロック単位で動作する CGRA の論理 4 列に写像された実際のデータフローを示している。

[0059] 図 15 は、実施例 1 と実施例 2 の演算ユニットにおけるアドレス計算部分の共通性を示すべく、実際に C 言語を用いたプログラムコードを示している。疎行列積の演算の場合 (図 15 における A のコード部分) も、マージソートの演算の場合 (図 15 における B のコード部分) も、プログラムコードの表現は同一であり、大小比較の部分の演算において、疎行列積の演算とマージソートの演算の何れの演算においてもアドレス計算部分の共通性があり、比較結果に基づいてアドレスを更新していることがわかる。

[0060] (その他の実施例)

(1) 本発明の疎行列積の演算方法の一実施形態について、図 16 の演算フローを参照して説明する。使用する演算ユニットは、3 入力パイプライン浮動小数点積和演算器とローカルメモリから構成され、積和演算器は、インデックスレジスタとデータレジスタを組にした第 1 の入力レジスタ及び第 2 の入力レジスタ、データレジスタである第 3 の入力レジスタ及び出力レジスタ

を備える。

疎行列積の演算フローは、まず、積和演算器の第1及び第2の入力レジスタの各インデックスレジスタにインデックス値Aとインデックス値Bを格納し、各データレジスタにデータAとデータBを格納する（ステップS01）。インデックス値Aは疎行列の列番号、インデックス値Bは疎行列の転置行列の列番号であり、各疎行列におけるゼロ値の要素を圧縮する（ステップS02）。各疎行列における非ゼロ値の要素がインデックス値とデータを組としてローカルメモリに記憶する（ステップS03）。次に、インデックス値Aとインデックス値Bを比較する（ステップS04）。

インデックス値A、Bが一致する場合（ $A=B$ ）、データAとデータBの乗算値と第3の入力レジスタのデータ値が加算されて積和演算器の出力レジスタに格納されると共に、第3の入力レジスタに戻す（ステップS06）。

インデックス値Aがインデックス値Bより小さいか等しい場合（ $A \leq B$ ）、第1のアドレス計算機構が、入力ベースアドレスレジスタに次要素参照に必要な値を加算する（ステップS07）。インデックス値Aがインデックス値Bより大きいか等しい場合（ $A \geq B$ ）、第2のアドレス計算機構が、入力ベースアドレスレジスタに次要素参照に必要な値を加算する（ステップS08）。ステップS07とステップS08の場合に、その後、アドレス加算結果を入力ベースアドレスレジスタに戻す（ステップS09）。これらのステップは、疎行列積が完了するまで繰り返され演算が完了する。

[0061] (2) 本発明のマージソートの演算方法の一実施形態について、図17の演算フローを参照して説明する。使用する演算ユニットは、3入力パイプライン浮動小数点積和演算器とローカルメモリから構成され、積和演算器は、インデックスレジスタとデータレジスタを組にした第1の入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備える。

マージソートの演算フローは、まず、積和演算器の第1及び第2の入力レジスタの各インデックスレジスタにインデックス値Aとインデックス値Bを

格納し、各データレジスタにデータAとデータBを格納する（ステップS11）。次に、インデックス値Aとインデックス値Bを比較する（ステップS12）。そして、インデックス値Aとインデックス値Bのデータの読み出し結果の大小関係に従い、後続の演算ユニットにローカルメモリの2つの読み出しアドレスと2つの読み出しデータを送り、アドレスとデータの各々の大小関係に従い、何れかのデータをローカルメモリにストアする（ステップS13）。これにより、マージソート全体の $\log N$ 段のうち、1段分のソート結果をローカルメモリにストアされる（ステップS14）。ローカルメモリのストア先アドレスを単調増加し（ステップS15）、ステップS11～S15の処理が、マージソート全体の $\log N$ 段まで繰り返され、マージソート演算が完了する。

### 産業上の利用可能性

[0062] 本発明は、メモリ容量に制約がある超小型AIアクセラレータに有用である。

### 符号の説明

- [0063]
- 1 第1のアドレス計算機構
  - 2 第2のアドレス計算機構
  - 3 演算ユニット
  - 5 外部メモリ
  - 31, 32, 33 入力レジスタ
  - 34 出力レジスタ
  - 15, 25, 35 ローカルメモリ

## 請求の範囲

[請求項1] 3入力パイプライン浮動小数点積和演算器とローカルメモリから構成される演算ユニットであって、

積和演算器は、インデックスレジスタとデータレジスタを組にした第1の入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備え、

積和演算器の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納され、

インデックス値Aとインデックス値Bが比較され、値が一致する場合には、データAとデータBの乗算値と第3の入力レジスタのデータ値が加算されて積和演算器の出力レジスタに格納されると共に、第3の入力レジスタに戻ることを特徴とする演算ユニット。

[請求項2] 疎行列Aと疎行列Bの行列積の演算において、インデックス値Aは疎行列の列番号で、インデックス値Bは疎行列の転置行列の列番号であり、

各疎行列におけるゼロ値の要素が圧縮され、非ゼロ値の要素がインデックス値とデータを組としてローカルメモリに記憶される請求項1に記載の演算ユニット。

[請求項3] クロック信号のタイミング $n+1$ で、積和演算器の第1段の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納され、

クロック信号のタイミング $n+2$ で、インデックス値Aとインデックス値Bが比較され、値が異なる場合には定数Cが積和演算器の第1段の出力データレジスタに格納され、値が一致する場合にはデータA

とデータ B の乗算が実行されて積和演算器の第 1 段の出力データレジスタに格納され、

クロック信号のタイミング  $n + 3$  で、積和演算器の第 2 段の出力データレジスタに加算結果が伝搬され、

クロック信号のタイミング  $n + 4$  で、積和演算器の第 3 段の出力データレジスタに正規化済の積和演算結果が伝搬され、

クロック信号のタイミング  $n + 5$  で、積和演算器の第 3 段の出力データレジスタから積和演算器の第 1 段の第 3 の入力レジスタであるデータレジスタに戻される、ことを特徴とする請求項 1 又は 2 に記載の演算ユニット。

[請求項4]

入力ベースアドレスレジスタと入力オフセットアドレスレジスタを含む第 1 のアドレス計算機構および第 2 のアドレス計算機構を備え、

第 1 および第 2 のアドレス計算機構は、各々の入力ベースアドレスレジスタおよび入力オフセットアドレスレジスタにアドレス情報をロードし、ローカルメモリから読み出したインデックス値 A とインデックス値 B を比較し、第 1 のアドレス計算機構は、 $A \leq B$  の場合に前記入力ベースアドレスレジスタに次要素参照に必要な値を加算し、第 2 のアドレス計算機構は、 $A \geq B$  の場合に前記入力ベースアドレスレジスタに次要素参照に必要な値を加算し、アドレス加算結果を前記入力ベースアドレスレジスタに戻すことを特徴とする請求項 1 ~ 3 の何れかに記載の演算ユニット。

[請求項5]

クロック信号の第 1 のタイミングで、各々の入力ベースアドレスレジスタおよび入力オフセットアドレスレジスタにアドレス情報をロードし、

クロック信号の第 2 のタイミングで、ローカルメモリから読み出したインデックス値 A とインデックス値 B を比較し、第 1 のアドレス計算機構は、 $A \leq B$  の場合に前記入力ベースアドレスレジスタに次要素参照に必要な値を加算し、第 2 のアドレス計算機構は、 $A \geq B$  の場合

に前記入力ベースアドレスレジスタに次要素参照に必要な値を加算し、

クロック信号の第3および第4のタイミングで、アドレス加算結果をデータ保持レジスタに伝搬させ、

クロック信号の次のタイミングで、前記データ保持レジスタの内容を前記入力ベースアドレスレジスタに戻すことを特徴とする請求項4に記載の演算ユニット。

[請求項6]

第1及び第2のアドレス計算機構は、ローカルメモリを複数空間に分割するアドレスマスク機構を備え、

クロック信号の第1のタイミングで、前記入力オフセットアドレスレジスタに値をセットし、

クロック信号の第2のタイミングで、前記入力オフセットアドレスレジスタの値をデータ保持レジスタに伝搬させ、

クロック信号の第3のタイミングで、各々の前記アドレス加算結果と前記入力オフセットアドレスレジスタの値を加算し、

クロック信号の第4のタイミングで、前記アドレスマスク機構により互いに異なる空間に分離してローカルメモリ参照のためのアドレスラッチに格納し、

クロック信号の第5のタイミングで、ローカルメモリから読み出したインデックスとデータの組をメモリ出力レジスタに格納し、

クロック信号の次のタイミングで、前記メモリ出力レジスタから読み出し、前記インデックス値として利用することを特徴とする請求項5の演算ユニット。

[請求項7]

3入力パイプライン浮動小数点積和演算器とローカルメモリから構成される演算ユニットを用いる演算方法であって、

インデックスレジスタとデータレジスタを組にした第1の入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備える積和演算器を用い、

積和演算器の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納されるステップと、インデックス値Aとインデックス値Bが比較されるステップと、インデックス値Aとインデックス値Bが一致する場合には、データAとデータBの乗算値と第3の入力レジスタのデータ値が加算されて積和演算器の出力レジスタに格納されると共に、第3の入力レジスタに戻すステップ、を備えることを特徴とする演算方法。

[請求項8] 前記演算方法が、疎行列Aと疎行列Bの行列積の演算方法であって、

インデックス値Aは疎行列の列番号で、インデックス値Bは疎行列の転置行列の列番号で、各疎行列におけるゼロ値の要素が圧縮されるステップと、

各疎行列における非ゼロ値の要素がインデックス値とデータを組としてローカルメモリに記憶されるステップを更に備えることを特徴とする請求項7に記載の演算方法。

[請求項9] 各々の入力ベースアドレスレジスタおよび入力オフセットアドレスレジスタにアドレス情報をロードするステップと、

ローカルメモリから読み出したインデックス値Aとインデックス値Bを比較するステップと、

第1のアドレス計算機構が、 $A \leq B$ の場合に前記入力ベースアドレスレジスタに次要素参照に必要な値を加算するステップ、或いは、第2のアドレス計算機構が、 $A \geq B$ の場合に前記入力ベースアドレスレジスタに次要素参照に必要な値を加算するステップと、

アドレス加算結果を前記入力ベースアドレスレジスタに戻すステップを更に備えることを特徴とする請求項7又は8に記載の演算方法。

[請求項10] 3入力パイプライン浮動小数点積和演算器とローカルメモリから構

成される演算ユニットを用いるマージソートの演算方法であって、

インデックスレジスタとデータレジスタを組にした第1の入力レジスタ及び第2の入力レジスタ、データレジスタである第3の入力レジスタ及び出力レジスタを備える積和演算器を用い、

積和演算器の第1の入力レジスタ及び第2の入力レジスタには、ローカルメモリから読み出したインデックス値Aとインデックス値Bが各インデックスレジスタに格納され、ローカルメモリから読み出したデータAとデータBが各データレジスタに格納されるステップと、

インデックス値Aとインデックス値Bが比較されるステップと、

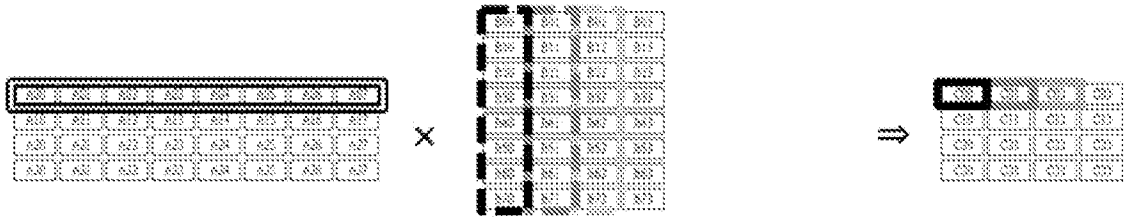
インデックス値Aとインデックス値Bのデータの読み出し結果の大小関係に従い、後続の演算ユニットにローカルメモリの2つの読み出しアドレスと2つの読み出しデータを送り、アドレスとデータの各々の大小関係に従い、何れかのデータをストアするステップと、

マージソート全体の $\log N$ 段のうち、1段分のソート結果をローカルメモリにストアするステップと、

ローカルメモリのストア先アドレスを単調に増加するステップ、を備えることを特徴とする演算方法。

[図1]

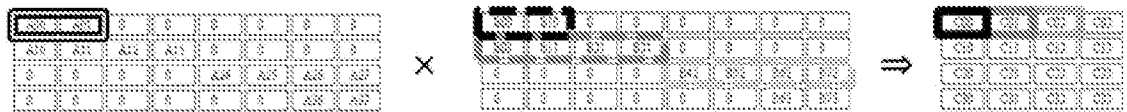
(1)  $A * B \Rightarrow C$



(2)  $A * B^T \Rightarrow C$



(3)  $A * B^T$  (疎行列)  $\Rightarrow C$



(4)  $A * B^T$  (圧縮表現)  $\Rightarrow C$



[図2]

疎行列A

疎行列B

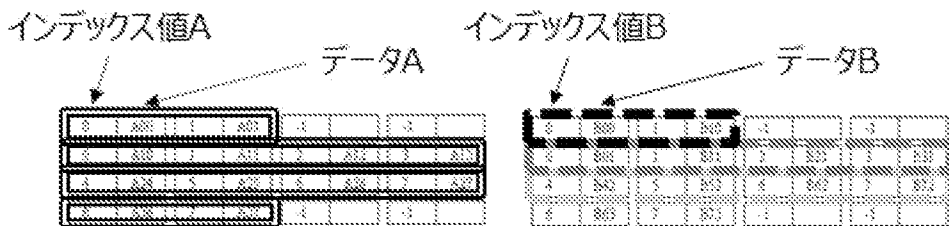
(1)



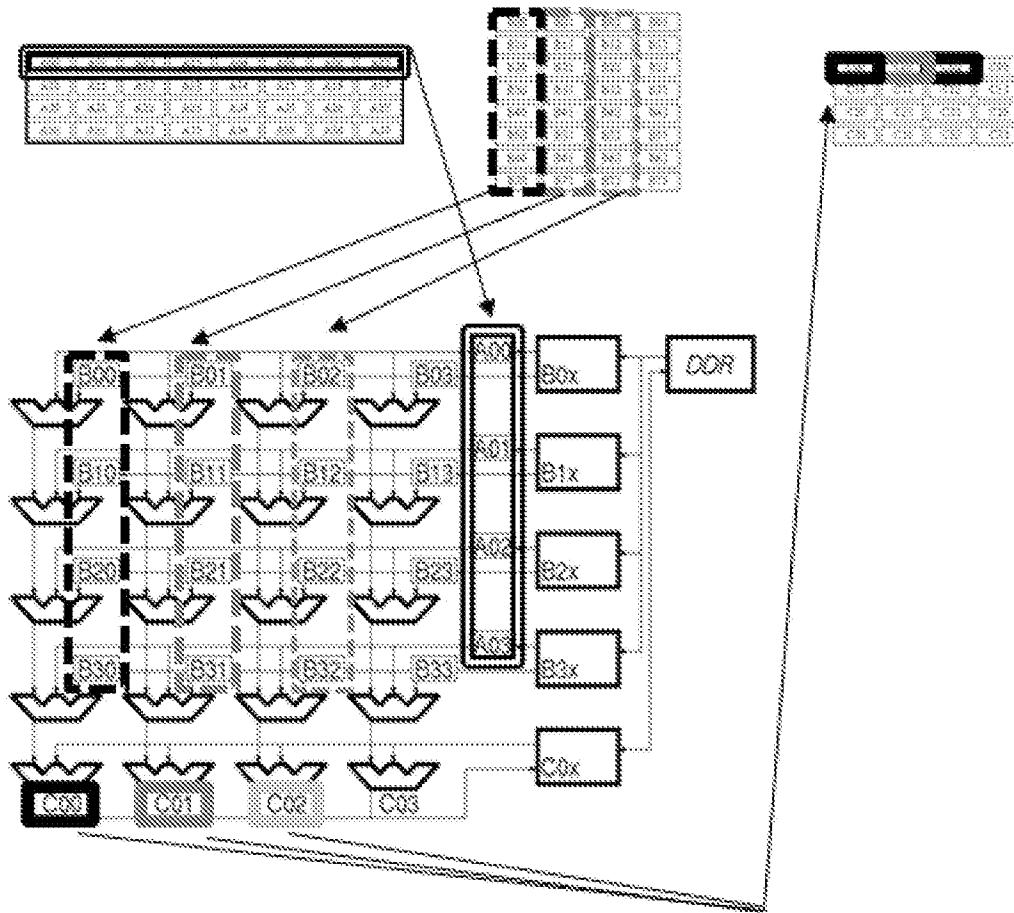
インデックス値を追加し、  
ゼロ値を圧縮

インデックス値を追加し、  
ゼロ値を圧縮

(2)

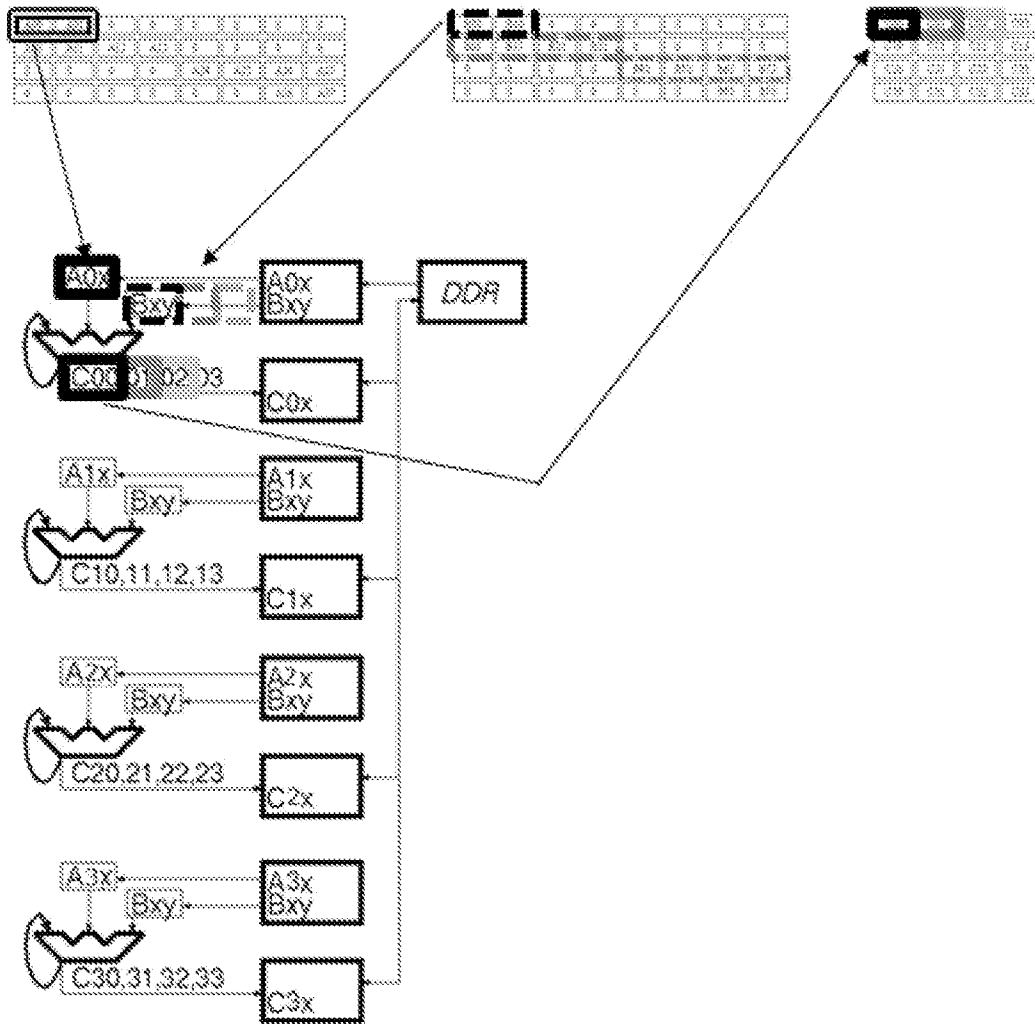


[図3]



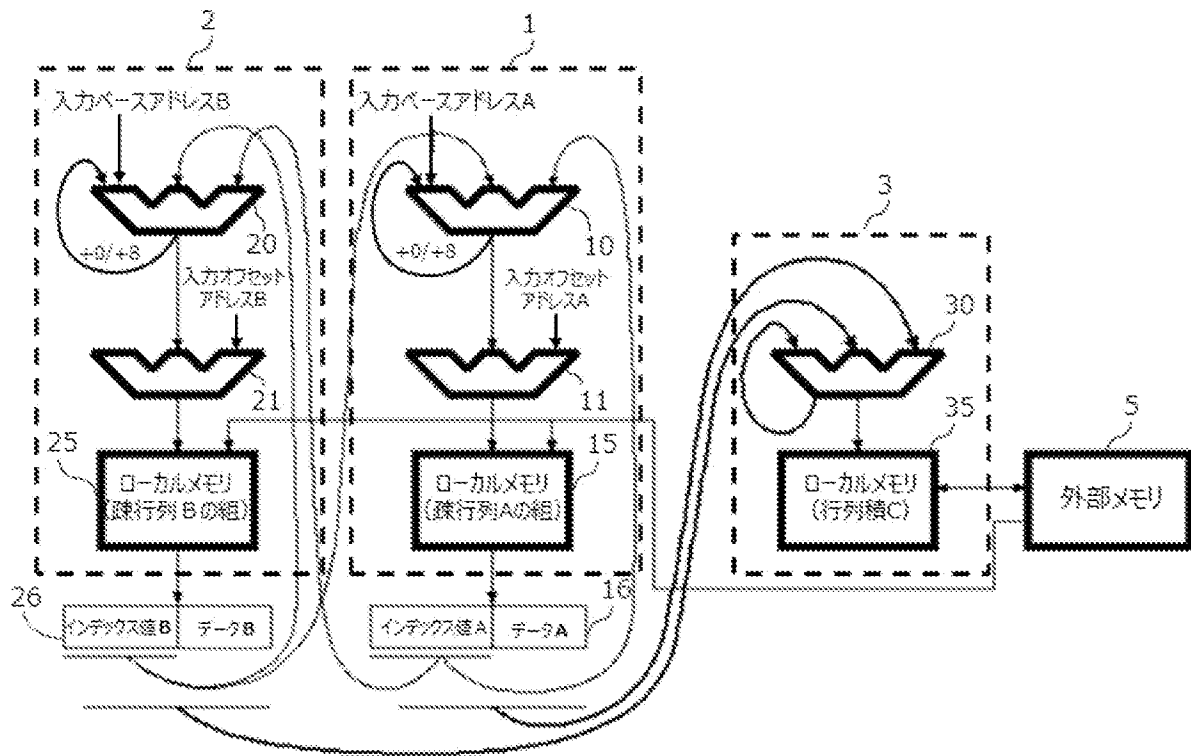
従来のCGRA演算の構成

[図4]

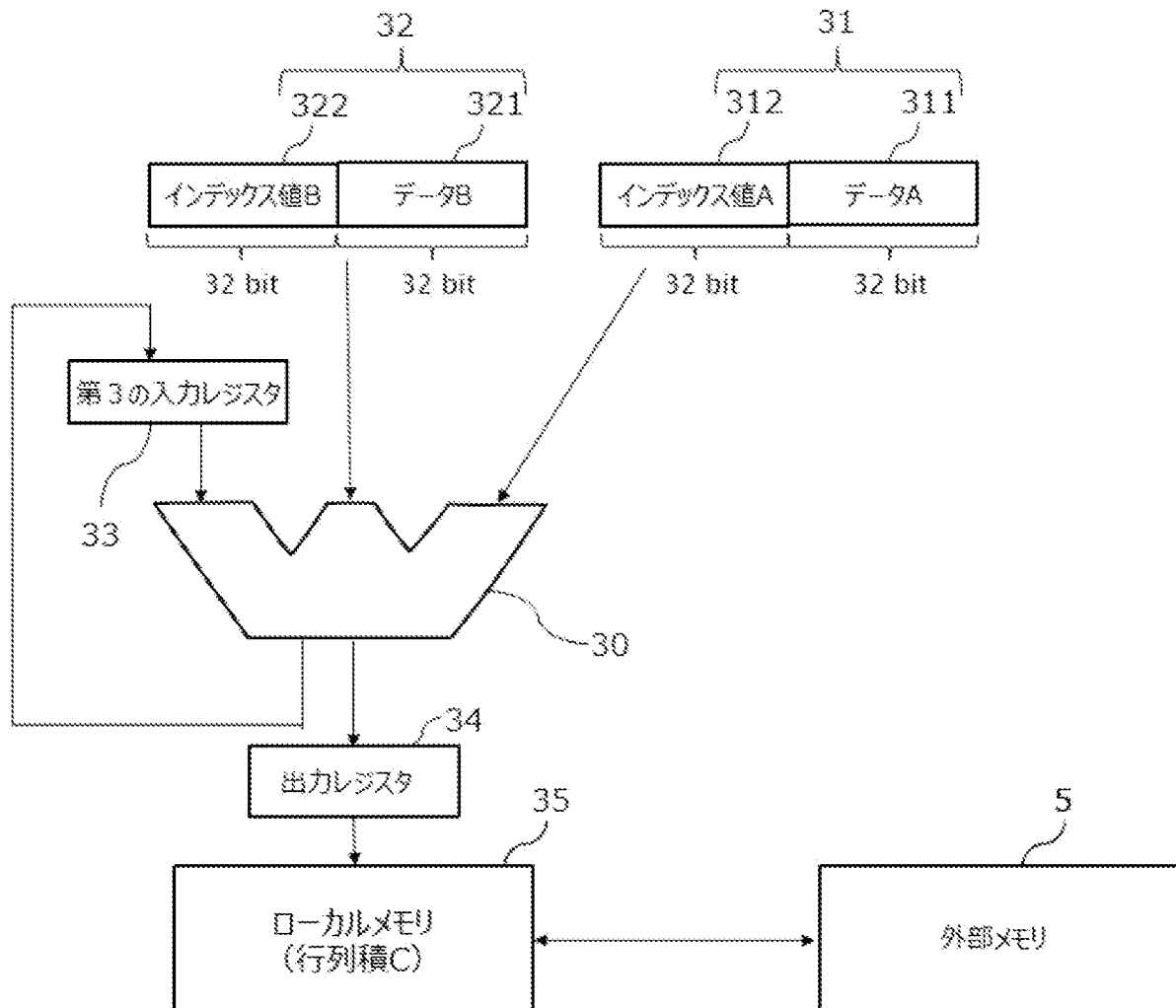


本発明のCGRA演算の構成

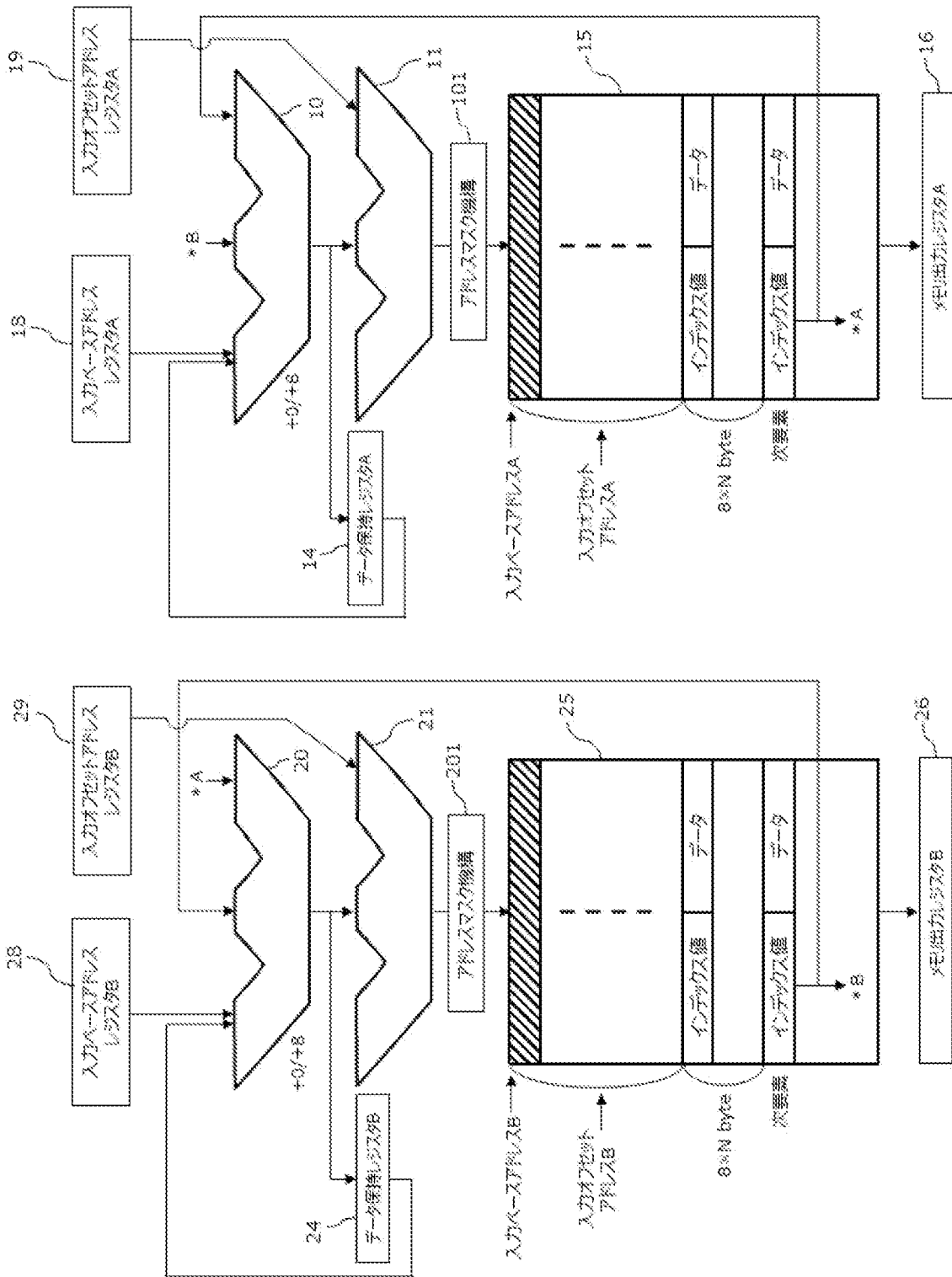
[図5]



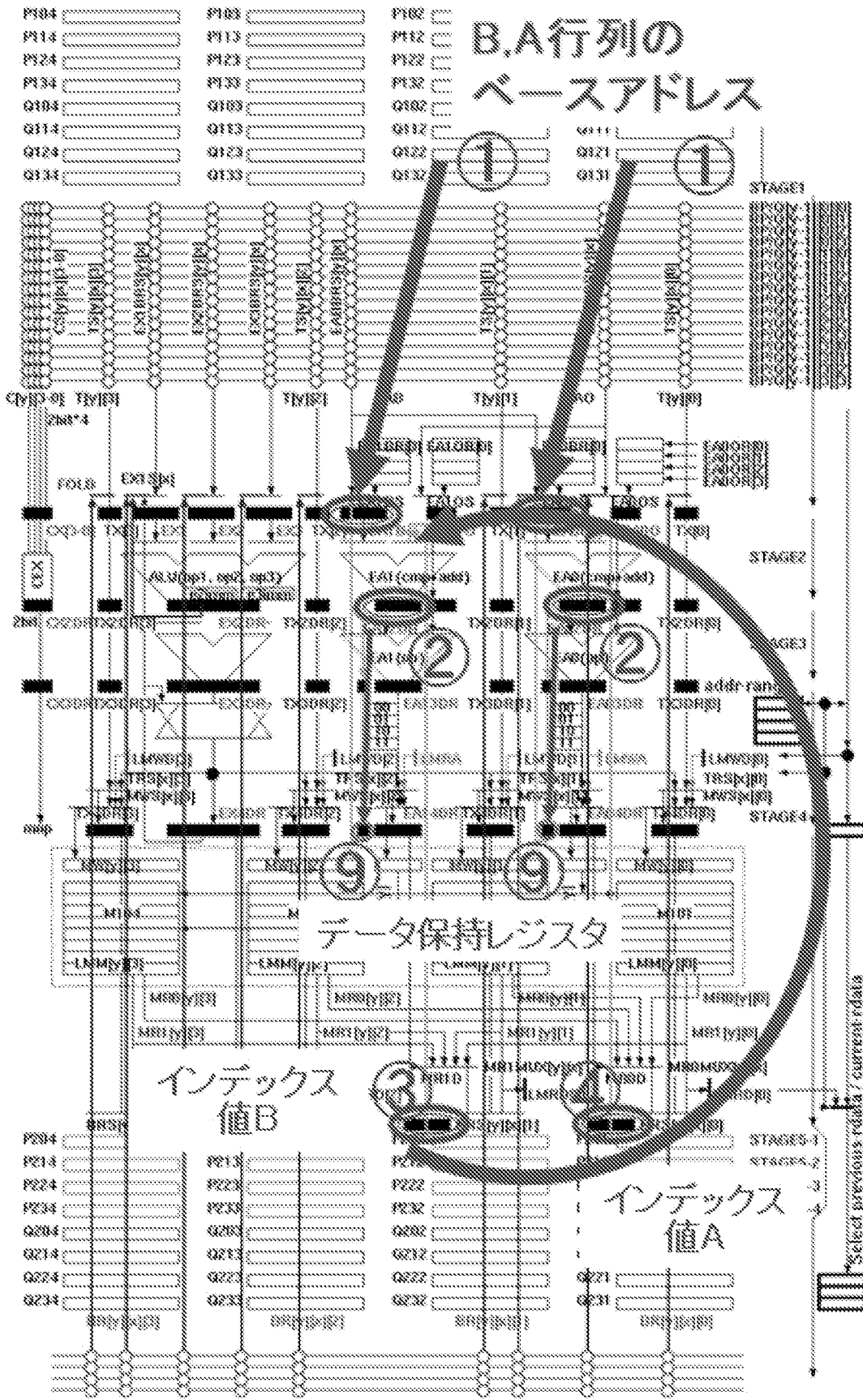
[図6]



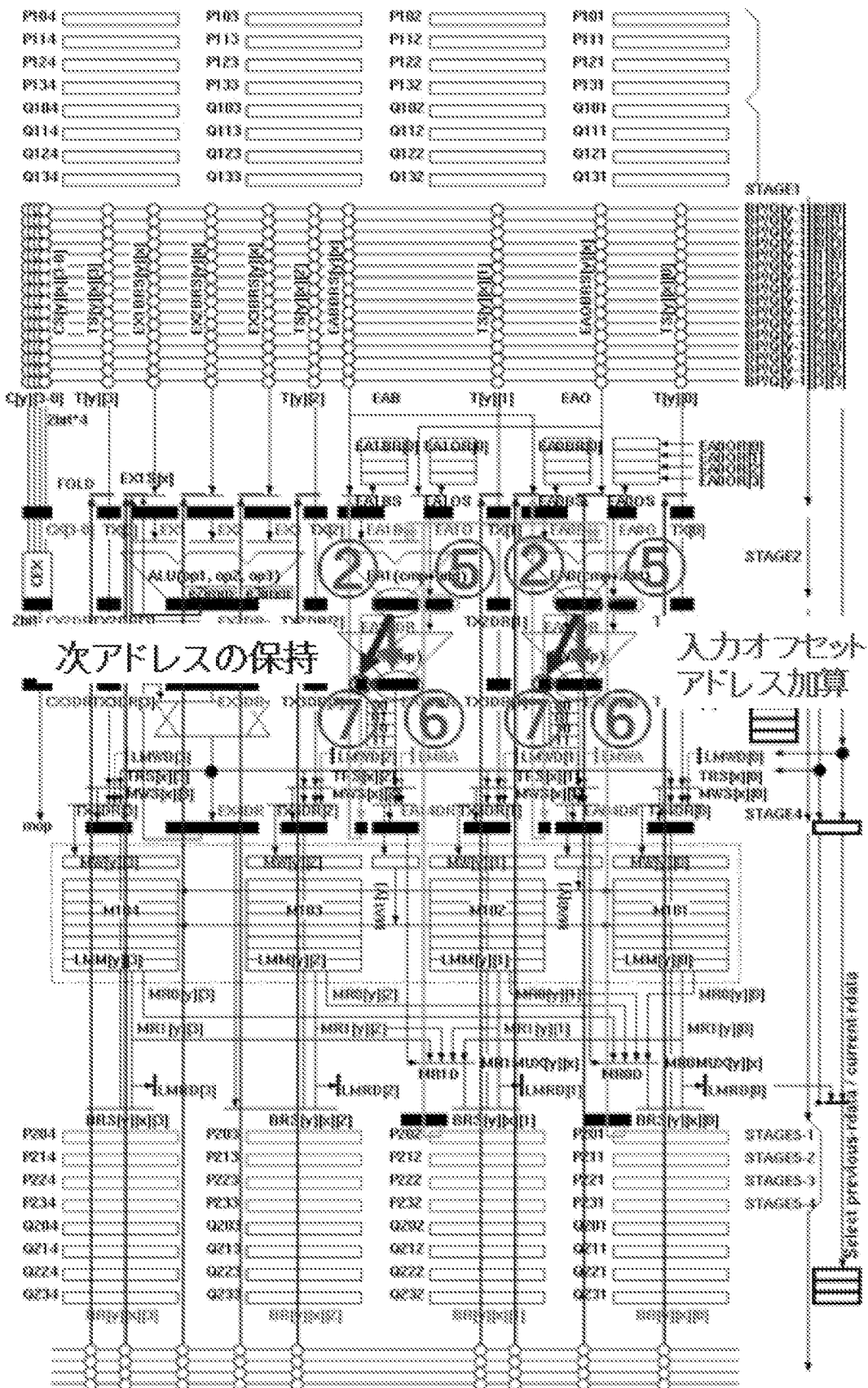
[図7]



[図8]

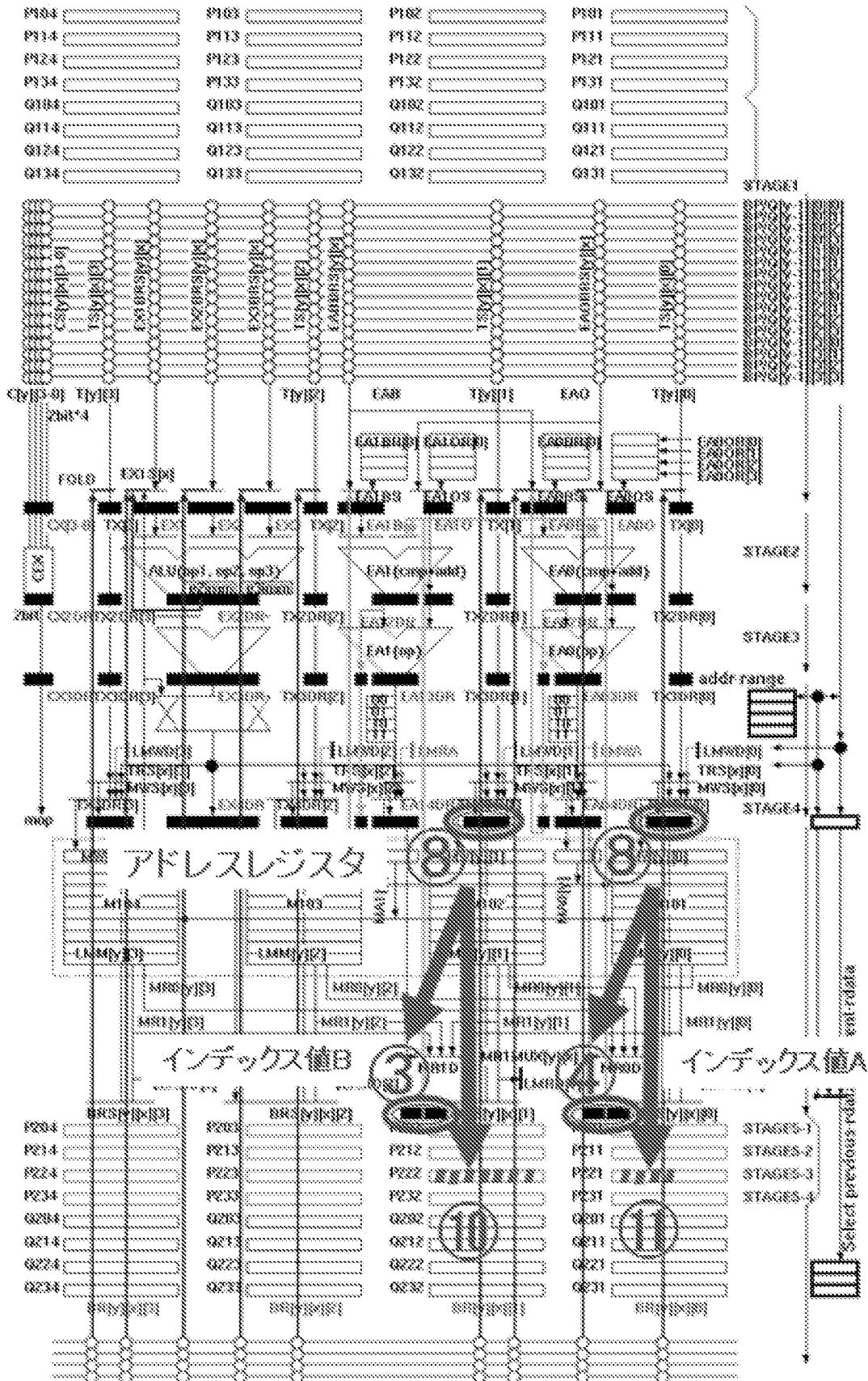


[図9]

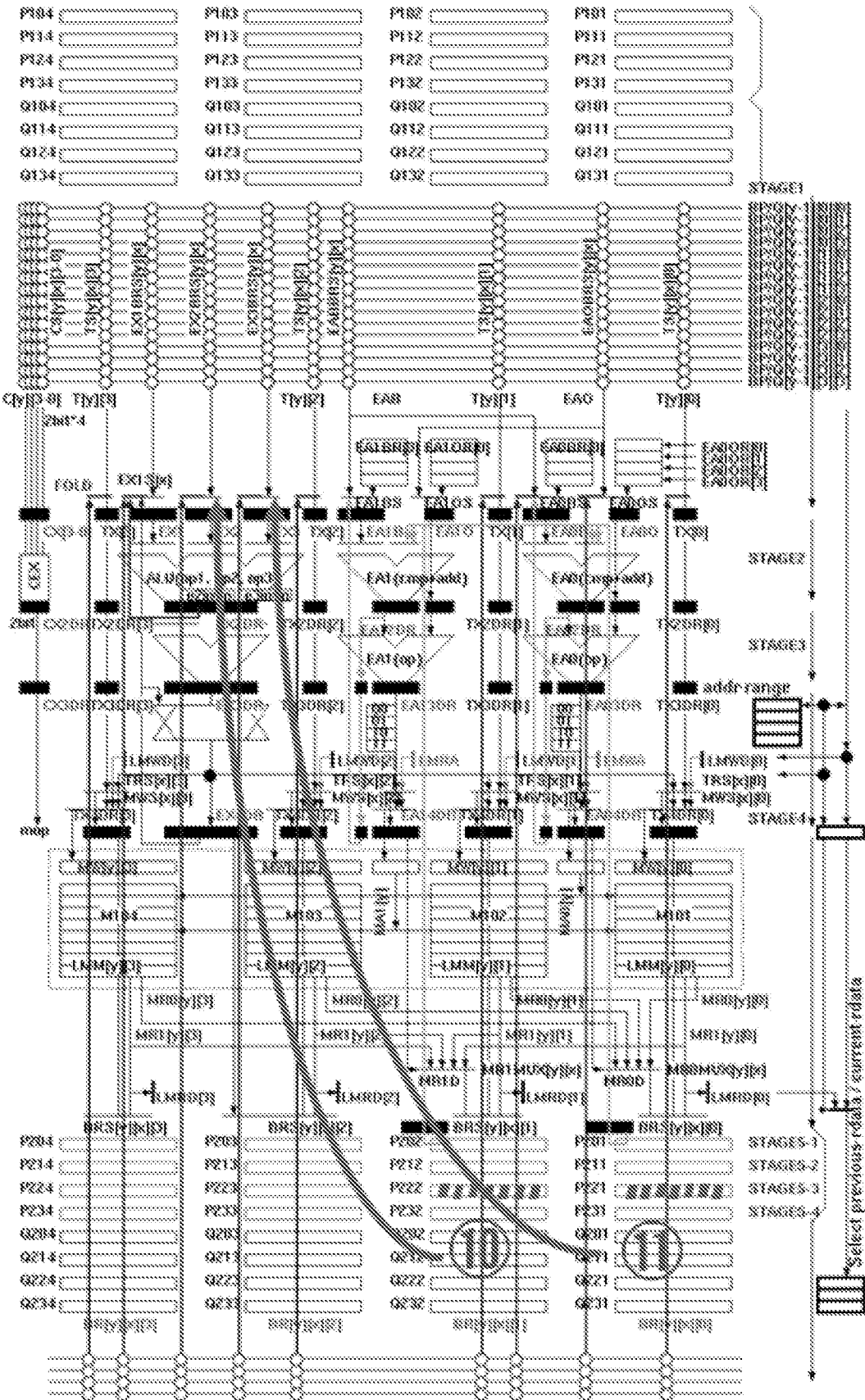




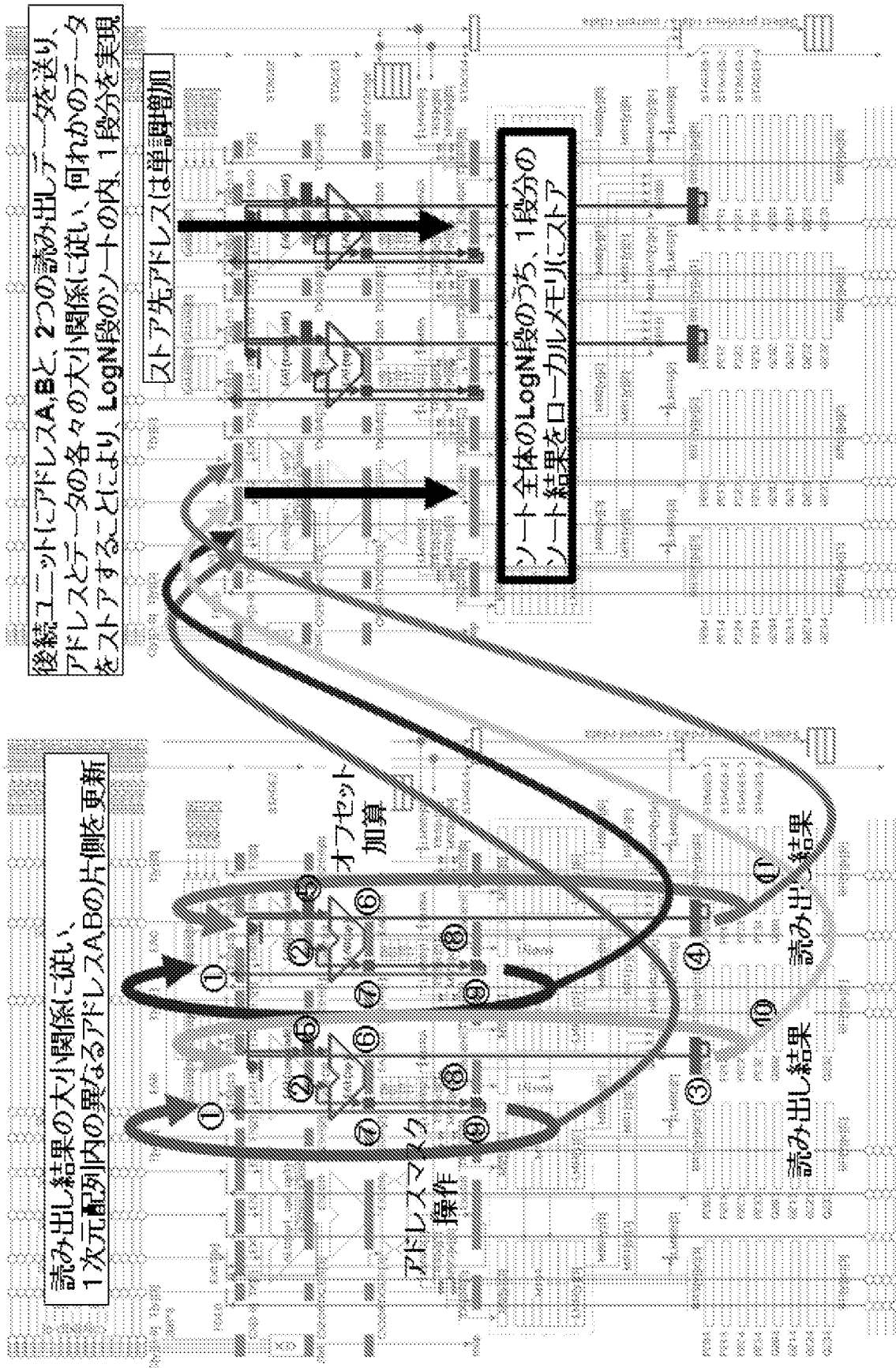
[図11]



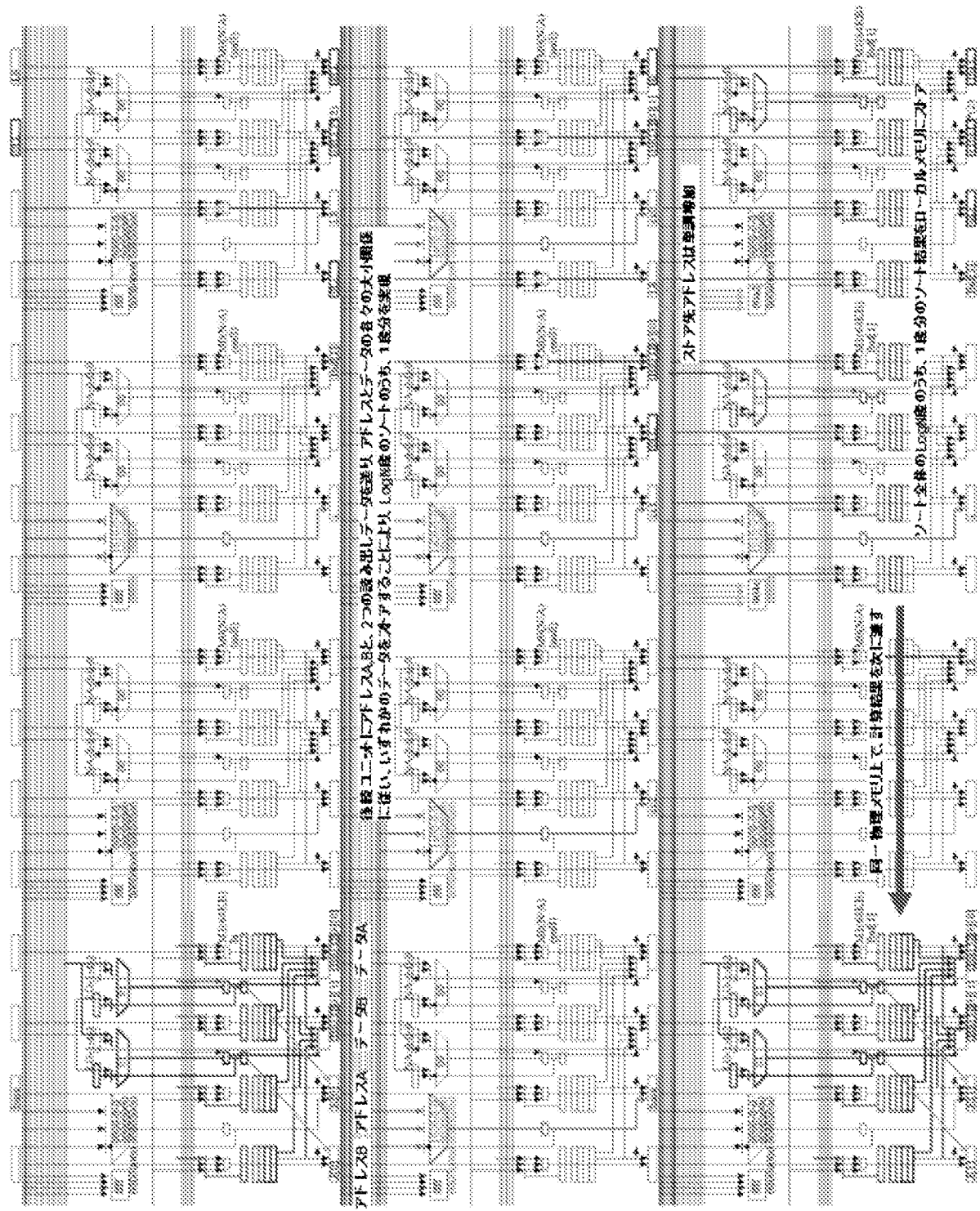
[図12]



[図13]



[図14]



[図15]

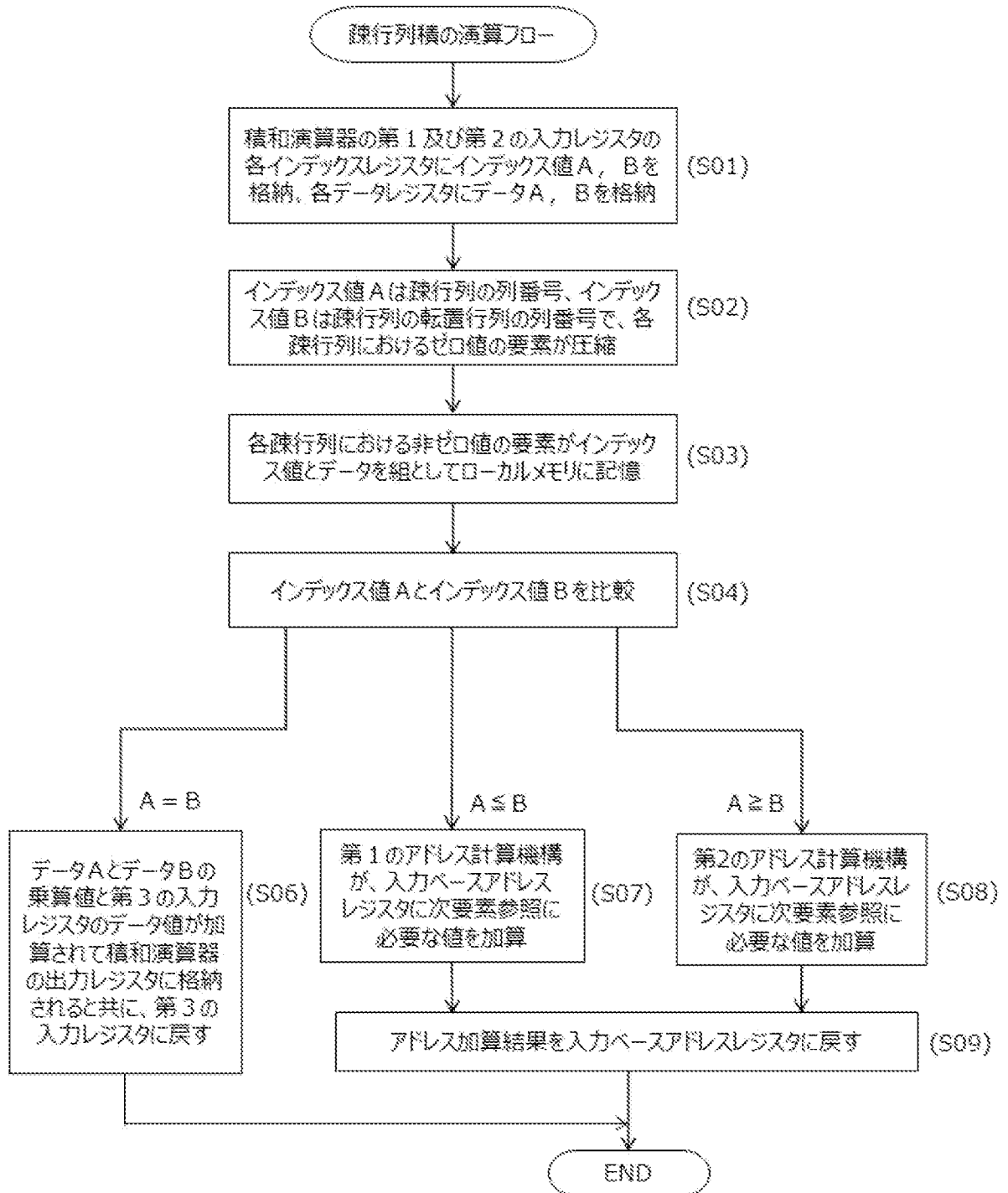
```

mex(uint op_mx, Uchar **d2, Uchar *base2, Uchar **d1, Uchar *base1, Ull limit, Ull ofs, Ull s2, Ull s1)
if (!limit) { /* sparse matrix */
    *d2 = base2 + ((s2!=0xfffffff && s2<=s1) ? ofs:0);
    *d1 = base1 + ((s1!=0xfffffff && s2>=s1) ? ofs:0);
}
else { /* merge sort */
    if ((base2==limit && base1+ofs==limit) || (base2+ofs==limit && base1==limit2)) {
        *d2 = 0;
        *d1 = limit;
    }
    else {
        *d2 = base2 + (base2!=limit && ((base1!=limit) || (base1==limit) ? ofs:0));
        *d1 = base1 + (base1!=limit2 && ((base2!=limit) || (base2==limit) ? ofs:0));
    }
}
}

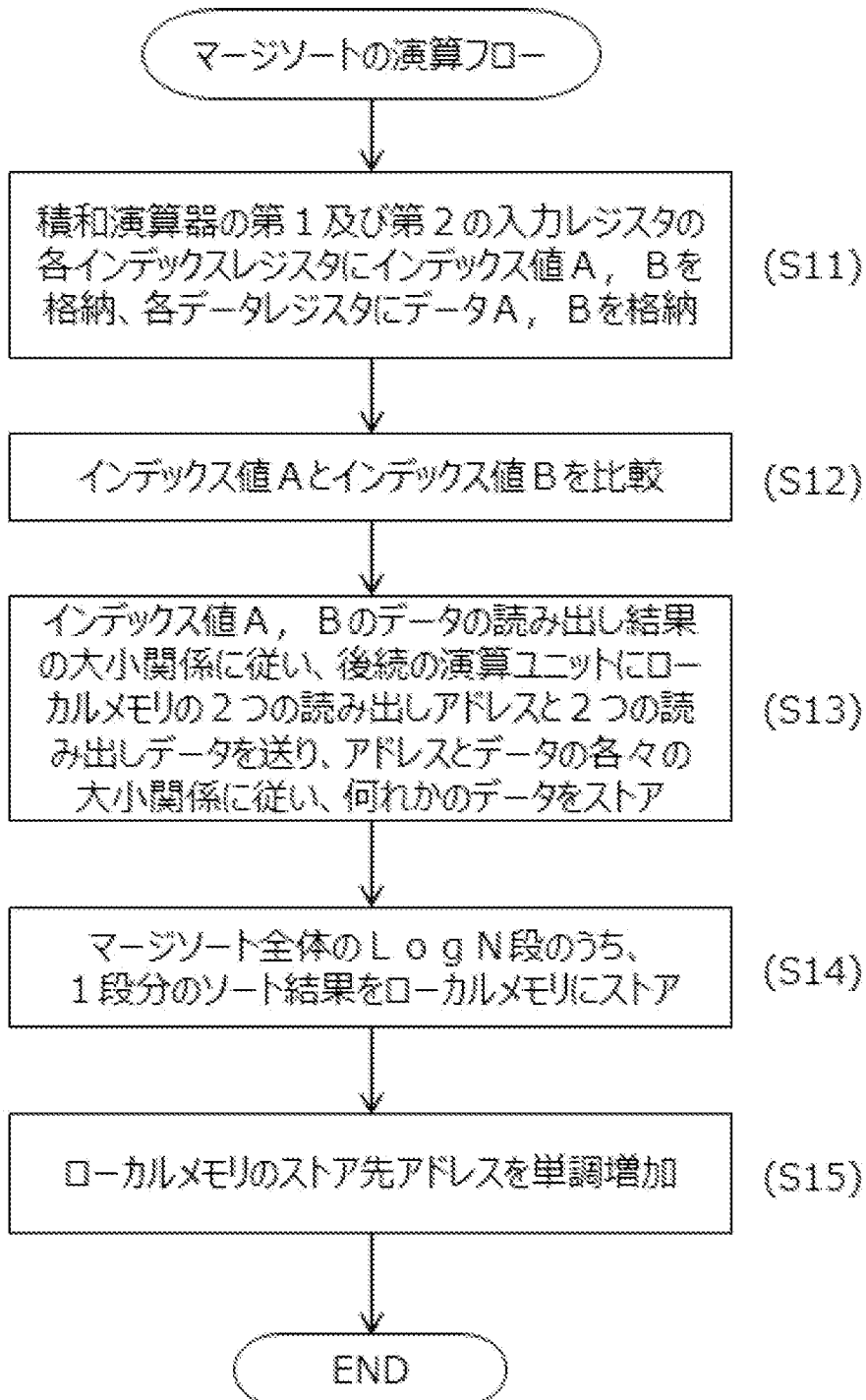
```

大小比較の部分の演算において、疎行列積とマージソートの何れもアドレス計算部分の共通性があり、比較結果に基づいてアドレスを更新

[図16]



[図17]



## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2022/046353

<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
<i>G06F 17/16</i> (2006.01)i; <i>G06F 7/24</i> (2006.01)i; <i>G06F 17/10</i> (2006.01)i FI: G06F17/16 P; G06F17/10 S; G06F7/24 Z		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) G06F17/16; G06F7/24; G06F17/10		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Published examined utility model applications of Japan 1922-1996 Published unexamined utility model applications of Japan 1971-2023 Registered utility model specifications of Japan 1996-2023 Published registered utility model applications of Japan 1994-2023		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	US 2019/0042538 A1 (INTEL CORPORATION) 07 February 2019 (2019-02-07) paragraphs [0043]-[0056], fig. 3-4B	1-2, 7-8
A	paragraphs [0043]-[0056], fig. 3-4B	3-6, 9-10
Y	JP 2008-234076 A (FUJITSU LIMITED) 02 October 2008 (2008-10-02) paragraphs [0006]-[0008], fig. 8	1-2, 7-8
A	paragraphs [0006]-[0008], fig. 8	3-6, 9-10
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search <b>06 February 2023</b>		Date of mailing of the international search report <b>14 February 2023</b>
Name and mailing address of the ISA/JP <b>Japan Patent Office (ISA/JP) 3-4-3 Kasumigaseki, Chiyoda-ku, Tokyo 100-8915 Japan</b>		Authorized officer  Telephone No.

**INTERNATIONAL SEARCH REPORT**  
**Information on patent family members**

International application No.

**PCT/JP2022/046353**

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
US	2019/0042538	A1	07 February 2019	(Family: none)	
JP	2008-234076	A	02 October 2008	US 2008/0228846 A1 paragraphs [0012]-[0021], fig. 8	

A. 発明の属する分野の分類（国際特許分類（IPC）） G06F 17/16(2006.01)i; G06F 7/24(2006.01)i; G06F 17/10(2006.01)i FI: G06F17/16 P; G06F17/10 S; G06F7/24 Z		
B. 調査を行った分野 調査を行った最小限資料（国際特許分類（IPC）） G06F17/16; G06F7/24; G06F17/10 最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2023年 日本国実用新案登録公報 1996-2023年 日本国登録実用新案公報 1994-2023年		
国際調査で使用した電子データベース（データベースの名称、調査に使用した用語）		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
Y	US 2019/0042538 A1 (INTEL CORPORATION) 07.02.2019 (2019-02-07) [0043]-[0056], FIGs.3-4B	1-2,7-8
A	[0043]-[0056], FIGs.3-4B	3-6,9-10
Y	JP 2008-234076 A (富士通株式会社) 02.10.2008 (2008-10-02) 【0006】～【0008】，【図8】	1-2,7-8
A	【0006】～【0008】，【図8】	3-6,9-10
<input type="checkbox"/> C欄の続きにも文献が列挙されている。 <input checked="" type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー “A” 特に関連のある文献ではなく、一般的な技術水準を示すもの “E” 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの “L” 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献（理由を付す） “O” 口頭による開示、使用、展示等に言及する文献 “P” 国際出願日前で、かつ優先権の主張の基礎となる出願の日の後に公表された文献 “T” 国際出願日又は優先日後に公表された文献であって出願と抵触するものではなく、発明の原理又は理論の理解のために引用するもの “X” 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの “Y” 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの “&” 同一パテントファミリー文献		
国際調査を完了した日	06.02.2023	国際調査報告の発送日 14.02.2023
名称及びあて先 日本国特許庁(ISA/JP) 〒100-8915 日本国 東京都千代田区霞が関三丁目4番3号	権限のある職員（特許庁審査官）  田中 幸雄 5B 9191  電話番号 03-3581-1101 内線 3545	

国際調査報告  
パテントファミリーに関する情報

国際出願番号

PCT/JP2022/046353

引用文献	公表日	パテントファミリー文献	公表日
US 2019/0042538 A1	07.02.2019	(ファミリーなし)	
JP 2008-234076 A	02.10.2008	US 2008/0228846 A1 [0012]-[0021], FIG. 8	