

(19)日本国特許庁(JP)

## (12)特許公報(B2)

(11)特許番号

特許第7187635号  
(P7187635)

(45)発行日 令和4年12月12日(2022.12.12)

(24)登録日 令和4年12月2日(2022.12.2)

(51)国際特許分類

F I

G 0 6 F 17/16 (2006.01)

G 0 6 F 17/16

P

G 0 6 F 17/16

Q

G 0 6 F 17/16

A

請求項の数 15 外国語出願 (全21頁)

(21)出願番号	特願2021-147894(P2021-147894)	(73)特許権者	502208397
(22)出願日	令和3年9月10日(2021.9.10)		グーグル エルエルシー
(62)分割の表示	特願2019-85959(P2019-85959)の 分割		G o o g l e L L C
原出願日	平成28年12月26日(2016.12.26)		アメリカ合衆国 カリフォルニア州 9 4
(65)公開番号	特開2022-781(P2022-781A)		0 4 3 マウンテン ビュー アンフィシ
(43)公開日	令和4年1月4日(2022.1.4)		アター パークウェイ 1 6 0 0
審査請求日	令和3年10月8日(2021.10.8)		1 6 0 0 A m p h i t h e a t r e P
(31)優先権主張番号	15/016,420		a r k w a y 9 4 0 4 3 M o u n t a
(32)優先日	平成28年2月5日(2016.2.5)	(74)代理人	i n V i e w , C A U . S . A .
(33)優先権主張国・地域又は機関	米国(US)		110001195
			弁理士法人深見特許事務所
		(72)発明者	ラビ・ナラヤナスワミ
			アメリカ合衆国、9 4 0 4 3 カリフォ
			ルニア州、マウンテン・ビュー、アンフ
			ィシアター・パークウェイ、1 6 0 0
			最終頁に続く

(54)【発明の名称】 疎要素を密行列に変換するためのシステムおよび方法

## (57)【特許請求の範囲】

## 【請求項 1】

疎要素を密行列に変換するためのシステムであって、  
 複数の疎要素アクセスユニットを備え、前記複数の疎要素アクセスユニットの各疎要素  
 アクセスユニットは、  
 それぞれの制御信号を受信し、  
 前記それぞれの制御信号に基づいて、前記疎要素アクセスユニットに対応するデータ片  
 に格納された疎要素にアクセスし、  
 前記データ片から得られる前記疎要素に適用される変換に基づいて出力密行列を生成し、  
 前記出力密行列を前記システムのノードネットワークに与えるよう構成されている、シ  
 ステム。

10

## 【請求項 2】

前記それぞれの制御信号に対応する命令を受けよう構成される疎密変換ユニットをさ  
 らに備え、  
 前記複数の疎要素アクセスユニットは、前記疎密変換ユニットに位置する、請求項 1 に  
 記載のシステム。

## 【請求項 3】

前記疎密変換ユニットは、行次元および列次元を含む多次元疎密変換ユニットである、  
 請求項 2 に記載のシステム。

## 【請求項 4】

20

前記複数の疎要素アクセスユニットは、前記多次元疎密変換ユニットのそれぞれの次元に沿って配置されている、請求項 3 に記載のシステム。

【請求項 5】

前記複数の疎要素アクセスユニットの各々は、

前記疎要素に前記変換を適用して前記出力密行列を生成するように構成されたそれぞれの第 1 のユニットを含む、請求項 2 に記載のシステム。

【請求項 6】

前記第 1 のユニットは連結ユニットであり、

前記変換は連結演算に基づく、請求項 5 に記載のシステム。

【請求項 7】

前記第 1 のユニットは圧縮 / 伸長ユニットであり、

前記変換は、前記疎要素を圧縮する演算に基づく、請求項 5 に記載のシステム。

【請求項 8】

複数の疎要素アクセスユニットを含むシステムを用いて疎要素を密行列に変換するための方法であって、

疎要素アクセスユニットがそれぞれの制御信号を受信することと、

前記それぞれの制御信号に基づいて、前記疎要素アクセスユニットが、前記疎要素アクセスユニットに対応するデータ片に格納された疎要素にアクセスすることと、

前記データ片から得られる前記疎要素に前記疎要素アクセスユニットによって適用される変換に基づいて、出力密行列を生成することと、

前記出力密行列を前記システムのノードネットワークに与えることとを含む、方法。

【請求項 9】

前記システムは、命令を受けるように構成された疎密変換ユニットを含み、

前記複数の疎要素アクセスユニットは、前記疎密変換ユニットに位置しており、

前記方法は、前記疎密変換ユニットが、前記複数の疎要素アクセスユニットの各々についてそれぞれの制御信号に対応する命令を受けることを含む、請求項 8 に記載の方法。

【請求項 10】

前記疎密変換ユニットは、行次元および列次元を含む多次元疎密変換ユニットである、請求項 9 に記載の方法。

【請求項 11】

前記複数の疎要素アクセスユニットは、前記多次元疎密変換ユニットのそれぞれの次元に沿って配置される、請求項 10 に記載の方法。

【請求項 12】

前記複数の疎要素アクセスユニットの各々は、前記疎要素を変換するように構成されたそれぞれの第 1 のユニットを含み、

前記方法は、前記第 1 のユニットが前記疎要素に前記変換を適用して前記出力密行列を生成することを含む、請求項 9 に記載の方法。

【請求項 13】

前記第 1 のユニットは連結ユニットであり、

前記変換を適用することは、

前記疎要素に連結演算を適用することと、

前記連結演算に基づいて前記疎要素を連結して前記出力密行列を生成することとを含む、請求項 12 に記載の方法。

【請求項 14】

前記第 1 のユニットは圧縮 / 伸長ユニットであり、

前記変換を適用することは、

前記疎要素に圧縮演算を適用することと、

前記圧縮演算に基づいて前記疎要素を圧縮して前記出力密行列を生成することとを含む、請求項 12 に記載の方法。

【請求項 15】

システムの疎要素アクセスユニットによって疎要素を密行列に変換するために用いられる命令を格納したコンピュータプログラムであって、前記命令は、請求項 8 ～ 14 のいずれか 1 項に記載の方法を引き起こすようプロセッサによって実行可能である、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

発明の詳細な記載

背景

この明細書は、一般に、回路系を用いて行列を処理することに関する。

10

【発明の概要】

【発明が解決しようとする課題】

【0002】

概要

この明細書に記載される主題の 1 つの革新的な局面によれば、行列プロセッサを用いて、疎から密への、または密から疎への行列変換を実行することができる。一般に、高性能計算システムは、行列を処理するために線形代数ルーチンを用い得る。いくつかの例においては、行列のサイズは 1 つのデータストレージにはまるには大きすぎるかもしれず、行列の異なる部分は、分散型データストレージシステムの異なる位置に疎に格納され得る。行列をロードするために、計算システムの中央処理ユニットは、別の回路系に行列の異なる部分にアクセスするよう命令し得る。この回路系は、ネットワークポロジリーに従って構成された複数のメモリコントローラを含んでもよく、疎データは、予め定められるルールの組に基いて、区分され格納されてもよい。各メモリコントローラは、予め定められるルールの組に基いて疎データを集めて、疎データ上において同時計算を実行し、および、中央処理ユニットがその後の処理を実行するために、ともに連結することができる密行列を生成してもよい。

20

【課題を解決するための手段】

【0003】

一般に、この明細書に記載される主題の 1 つの革新的な局面は、疎要素を密行列に変換するためのシステムにおいて実施することができる。このシステムは、第 1 の密行列と関連付けられる疎要素をフェッチするよう構成された疎要素アクセスユニットの第 1 の群と、第 1 の密行列とは異なる第 2 の密行列と関連付けられる疎要素をフェッチするよう構成された疎要素アクセスユニットの第 2 の群とを備える。システムは、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を含む疎要素に基いた出力行列に対する要求を受取り、疎要素アクセスユニットの第 1 の群によってフェッチされる第 1 の密行列と関連付けられる疎要素を得、疎要素アクセスユニットの第 2 の群によってフェッチされる第 2 の密行列と関連付けられる疎要素を得、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を変換して、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を含む出力密行列を生成するよう構成される。

30

40

【0004】

これらおよび他の実現例は、各々、任意で以下の特徴の 1 つ以上を含むことができる。たとえば、疎要素アクセスユニットの第 1 の群は第 1 の疎要素アクセスユニットおよび第 2 の疎要素アクセスユニットを含んでもよい。第 1 の疎要素アクセスユニットは、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合をフェッチするよう構成されてもよい。第 2 の疎要素アクセスユニットは、第 1 の密行列と関連付けられる疎要素の第 2 の異なる部分集合をフェッチするよう構成されてもよい。

【0005】

第 1 の疎要素アクセスユニットは、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を含む複数個の疎要素に対する要求を受取り、要求を第 2

50

の疎要素アクセスユニットに送信するよう構成される。第 1 の疎要素アクセスユニットは、複数の疎要素のうちの特定の疎要素のアイデンティティが、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合のうちの 1 つのアイデンティティと一致する、と判断してもよい。複数の疎要素のうちの特定の疎要素のアイデンティティが、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合のうちの 1 つのアイデンティティと一致する、と判断することに応じて、第 1 の疎要素アクセスユニットは、特定の疎要素を含む第 1 の密行列と関連付けられる疎要素の第 1 の部分集合をフェッチするよう構成されてもよい。

【 0 0 0 6 】

第 1 の疎要素アクセスユニットは、第 1 のデータ片から第 1 の密行列と関連付けられる疎要素の第 1 の部分集合をフェッチするよう構成されてもよく、第 2 の疎要素アクセスユニットは、第 2 の異なるデータ片から第 1 の密行列と関連付けられる疎要素の第 2 の異なる部分集合をフェッチするよう構成されてもよい。第 1 の疎要素アクセスユニットは、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合を変換して第 3 の密行列を生成するよう構成されてもよく、第 2 の疎要素アクセスユニットは、第 3 の密行列を受け、第 2 の密行列と関連付けられる疎要素の第 2 の部分集合を変換して第 4 の密行列を生成し、第 3 の密行列を第 4 の密行列とともに変換して、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合および第 1 の密行列と関連付けられる疎要素の第 2 の部分集合を含む第 5 の密行列を生成するよう構成されてもよい。

【 0 0 0 7 】

疎要素アクセスユニットの第 1 の群および疎要素アクセスユニットの第 2 の群は二次元のメッシュ構成で配列されてもよい。疎要素アクセスユニットの第 1 の群および疎要素アクセスユニットの第 2 の群は二次元の円環面構成で配列されてもよい。第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素は多次元の行列であってもよく、出力密行列はベクトルであってもよい。

【 0 0 0 8 】

この明細書において記載される主題は、以下の利点の 1 つ以上を実現するように特定の実施の形態において実現することができる。ネットワークポロジーに従ってメモリコントローラユニットを接続することは、予め定められるルールの組に従う疎データの格納の区分化を可能にする。中央処理ユニットから別の回路系に疎密データロードタスクをシフトすることは、中央処理ユニットの計算帯域幅を増大し、システムの処理費を低減する。特殊化された回路系を用いることによって、疎データをフェッチするために密な線形代数に対して特殊化されるプロセッサの使用を回避することができる。分散型システムにおいて多数のメモリを同時に用いることによって、分散型システムにおいて利用可能な和集合帯域幅は、直列化を必要とし、集合の帯域幅上において単一のメモリキャップを有する単一のメモリバンクに対する帯域幅よりも高い。

【 0 0 0 9 】

この局面および他の局面の他の実現例は、計算機記憶装置上でエンコードされる、方法のアクションを実行するよう構成される、対応のシステム、装置およびコンピュータプログラムを含む。1 つ以上のコンピュータのシステムは、システムにインストールされ、動作でシステムにアクションを実行させるソフトウェア、ファームウェア、ハードウェアまたはそれらの組合せによってそのように構成することができる。1 つ以上のコンピュータプログラムは、データ処理装置によって実行されたとき、装置にアクションを実行させる命令を有することによって、そのように構成することができる。

【 0 0 1 0 】

この明細書に記載される主題の 1 つ以上の実現例の詳細は、添付の図面および以下の記載において述べられる。主題の他の潜在的な特徴、局面および利点は、記載、図面および特許請求の範囲から明らかになる。

【図面の簡単な説明】

【 0 0 1 1 】

【図 1】例示の計算システムのブロック図である。

10

20

30

40

50

- 【図 2 A】例示の疎密変換ユニットを示す図である。
- 【図 2 B】例示の疎密変換ユニットを示す図である。
- 【図 2 C】例示の疎密変換ユニットを示す図である。
- 【図 2 D】例示の疎密変換ユニットを示す図である。
- 【図 3 A】例示の疎要素アクセスユニットを示す図である。
- 【図 3 B】例示の疎要素アクセスユニットを示す図である。
- 【図 4】密行列を生成するためのプロセスの例を示すフローチャート図である。
- 【図 5】疎要素を密行列に変換するためのプロセスの例を示すフローチャート図である。
- 【発明を実施するための形態】

#### 【 0 0 1 2 】

さまざまな図面における同様の参照番号および指定は同様の要素を示す。

#### 詳細な記載

一般に、データは行列の形式において表すことができ、計算システムは線形代数アルゴリズムを用いてデータを操作し得る。行列は一次元のベクトルまたは多次元行列であり得る。行列は、データベーステーブルまたは変数などのようなデータ構造によって表されてもよい。しかしながら、行列のサイズが大きすぎると、1つのデータストレージに行列全体を格納することは可能ではないかもしれない。密行列は複数の疎要素に変換され得、各疎要素は異なるデータストレージに格納され得る。密行列の疎要素は行列であってもよく、行列のうちの小さな部分行列（たとえば単一値要素、行、列、または部分行列）のみが非零値を有する。計算システムが密行列にアクセスすることを必要とするときに、中央処理ユニット（CPU）は、データストレージの各々に到達するスレッドを開始して、格納された疎要素をフェッチしてもよく、そして、疎密変換を適用して密行列を戻す。しかしながら、それが疎要素すべてをフェッチするのにかかる時間の量は長いかもしれず、CPUの計算帯域幅は結果として十分に利用されないかもしれない。いくつかの場合では、計算システムはいくつかの密行列の疎要素にアクセスして新たな密行列を形成する必要があるかもしれず、それらの密行列は等しい次元を有さないかもしれない。異なる密行列の疎要素をフェッチするようデータストレージの各々に到達するスレッドと関連付けられるCPUアイドル時間は、異なる待ち時間に遭遇し得、さらに、計算装置の性能に望ましくない態様で影響を与えるかもしれない。いくつかの場合では、計算システムはいくつかの密行列の疎要素にアクセスして新たな密行列を形成する必要があるかもしれず、それらの疎要素は等しい次元を有さないかもしれない。異なる密行列の疎要素をフェッチするようデータストレージの各々に到達するスレッドと関連付けられるCPUアイドル時間は、異なる待ち時間に遭遇し得、さらに、計算装置の性能に望ましくない態様で影響を与えるかもしれない。CPUから分離しているハードウェア疎密変換ユニットは、プロセッサの計算帯域幅をCPU動作から独立した疎要素の収集および疎要素の密行列への変換によって、増大させ得る。

#### 【 0 0 1 3 】

図 1 は、1つ以上の密行列から疎要素を変換して密行列を生成するための例示の計算システム 100 のブロック図を示す。計算システム 100 は、処理ユニット 102、疎密変換ユニット 104 およびデータ片 106 a ~ 106 k を含み、k は 1 以上の整数である。一般に、処理ユニット 102 は、目標密行列へのアクセスのための命令を処理し、目標密行列を生成するために疎密変換ユニット 104 に命令 110 を送信する。疎密変換ユニット 104 はデータ片 106 a ~ 106 k の 1 つ以上から対応する疎要素 108 a ~ 108 n にアクセスし、n は 1 つ以上の整数である。疎密変換ユニット 104 は対応する疎要素 108 a ~ 108 n を用いて目標密行列 112 を生成し、目標密行列 112 をその後の処理のために処理ユニット 102 に与える。たとえば、疎要素 108 a ~ 108 n は異なるサイズを有する二次元の行列であり得、疎密変換ユニット 104 は疎要素 108 a ~ 108 n の各々をベクトルに変換することおよび n 個のベクトルを単一のベクトルに連結することによって目標密行列 112 を生成し得る。

#### 【 0 0 1 4 】

いくつかの実現例では、処理ユニット102は、目標密行列の更新ための命令を処理し、更新された密行列を疎密変換ユニット104に送信してもよい。疎密変換ユニット104は、更新された密行列を対応する疎要素に変換し、したがって、データ片106a~106kに格納された1つ以上の疎要素を更新してもよい。

#### 【0015】

処理ユニット102は計算システム100内で実行のために命令を処理するように構成される。処理ユニット102は1つ以上のプロセッサを含んでもよい。いくつかの実現例では、処理ユニット102は疎密変換ユニット104によって生成された目標密行列112を処理するように構成される。他のいくつかの実現例では、処理ユニット102は、疎密変換ユニット104に目標密行列112を生成することを要求するよう構成されてもよく、別の処理ユニットが目標密行列112を処理するように構成されてもよい。データ片106a~106kは疎要素108a~108nを含むデータを格納する。いくつかの実現例では、データ片106a~106kは単数または複数の揮発性記憶装置であってもよい。他のいくつかの実現例では、データ片106a~106kは単数または複数の不揮発性記憶装置であってもよい。データ片106a~106kは、さらに、ストレージエリアネットワークまたは他の構成における装置などのような別のコンピュータ読取可能媒体の形式であってもよい。データ片106a~106kは電氣的接続、光接続または無線接続を用いて、疎密変換ユニット104に結合されてもよい。いくつかの実現例では、データ片106a~106kは疎密変換ユニット104の一部であってもよい。

#### 【0016】

疎密変換ユニット104は疎要素に基いて密行列を判断するように構成される。いくつかの実現例では、疎密変換ユニット104は密行列に基いて疎要素の位置を判断するように構成されてもよい。いくつかの実現例では、図2A~図2Dを参照してより詳細に以下に記載されるように、疎密変換ユニット104は複数の相互接続された疎要素アクセスユニットを含んでもよい。

#### 【0017】

図2Aは例示の疎密変換ユニット200を示す。疎密変換ユニット200は疎密変換ユニット104に対応してもよい。疎密変換ユニット200は、 $M \times N$ 個の疎要素アクセスユニット $X_{1,1} \sim X_{M,N}$ を含み、それらは物理的にまたは論理的にM個の行およびN個の列に配列され、MおよびNは1以上の整数である。いくつかの実現例では、疎密変換ユニット200は、データを処理するように構成されたさらなる回路系を含んでもよい。一般に、疎密変換ユニット200は、密行列に対する要求を受け、疎要素アクセスユニット $X_{1,1} \sim X_{M,N}$ によってアクセス可能な対応する疎要素に基いて密行列を判断するよう構成される。一般に、各疎要素アクセスユニットは、指定される疎要素の組にアクセスするよう構成され、図3A~図3Bを参照してより詳細に以下に記載される。いくつかの実現例では、疎要素アクセスユニットは単一命令・多重データ(SIMD)処理装置であってもよい。

#### 【0018】

いくつかの実現例では、疎要素アクセスユニット $X_{1,1} \sim X_{M,N}$ は、二次元のメッシュ構成に物理的にまたは論理的に配列されてもよい。たとえば疎要素アクセスユニット $X_{1,1}$ は、疎要素アクセスユニット $X_{1,2}$ および $X_{2,1}$ に直接結合される。別の例として、疎要素アクセスユニット $X_{2,2}$ は、疎要素アクセスユニット $X_{2,1}$ 、 $X_{3,1}$ 、 $X_{2,3}$ および $X_{1,2}$ に直接結合される。2つの疎要素アクセスユニット間の結合は、電氣的接続、光接続、無線接続または任意の他の好適な接続であってもよい。

#### 【0019】

他のいくつかの実現例では、疎要素アクセスユニット $X_{1,1} \sim X_{M,N}$ は、二次元の円環面構成に物理的にまたは論理的に配列されてもよい。たとえば疎要素アクセスユニット $X_{1,1}$ は、疎要素アクセスユニット $X_{1,2}$ 、 $X_{2,1}$ 、 $X_{1,N}$ および $X_{M,1}$ に直接結合される。別の例として、疎要素アクセスユニット $X_{M,N}$ は、疎要素アクセスユニット $X_{M,N-1}$ 、 $X_{M-1,N}$ 、 $X_{M,1}$ および $X_{1,N}$ に直接結合される。

## 【 0 0 2 0 】

いくつかの実現例では、疎密変換ユニット 2 0 0 は、予め定められる条件の組に従って密行列から変換される疎要素を区分するよう構成されてもよい。疎要素アクセスユニット  $X_{1,1} \sim X_{M,N}$  の各行は、特定の密行列から変換される疎要素にアクセスするよう区分されてもよい。たとえば、疎密変換ユニット 2 0 0 は、コンピュータモデルの 1, 0 0 0 個の異なるデータベーステーブルに対応する密行列から変換される疎要素にアクセスするよう構成されてもよい。データベーステーブルの 1 つ以上は異なるサイズを有してもよい。疎要素アクセスユニットの第 1 番目の行 2 0 2 は、データベーステーブル 1 番 ~ データベーステーブル 1 0 0 番から変換される疎要素にアクセスするよう構成されてもよく、疎要素アクセスユニットの第 2 番目の行 2 0 4 は、データベーステーブル 1 0 1 番 ~ データベーステーブル 3 0 0 番から変換される疎要素にアクセスするよう構成されてもよく、疎要素アクセスユニットの M 番目の行 2 0 6 は、データベーステーブル 7 5 1 番 ~ データベーステーブル 1, 0 0 0 番から変換される疎要素にアクセスするよう構成されてもよい。いくつかの実現例では、区分は、疎密変換ユニット 2 0 0 を用いて、プロセッサが疎要素にアクセスする前に、ハードウェア命令によって構成されてもよい。

10

## 【 0 0 2 1 】

疎要素アクセスユニット  $X_{1,1} \sim X_{M,N}$  の各列は、特定の密行列から変換される疎要素の部分集合にアクセスするよう区分されてもよい。たとえば、データベーステーブル 1 番に対応する密行列は 1, 0 0 0 個の疎要素に変換されてもよく、1, 0 0 0 個の疎要素は上に記載されるような第 1 番目の行 2 0 2 によってアクセス可能である。疎要素アクセスユニット  $X_{1,1}$  は、データベーステーブル 1 番の疎要素 1 番 ~ 2 0 0 番にアクセスするよう構成されてもよく、疎要素アクセスユニット  $X_{1,2}$  は、データベーステーブル 1 番の疎要素 2 0 1 番 ~ 5 0 0 番にアクセスするよう構成されてもよい。別の例として、データベーステーブル 2 番に対応する密行列は 5 0 0 個の疎要素に変換されてもよく、5 0 0 個の疎要素は上に記載されるような第 1 番目の行 2 0 2 によってアクセス可能である。疎要素アクセスユニット  $X_{1,1}$  は、データベーステーブル 2 番の疎要素 1 番 ~ 5 0 番にアクセスするよう構成されてもよく、疎要素アクセスユニット  $X_{1,2}$  は、データベーステーブル 2 番の疎要素 5 1 番 ~ 2 0 0 番にアクセスするよう構成されてもよい。別の例として、データベーステーブル 1, 0 0 0 番に対応する密行列は 1 0, 0 0 0 個の疎要素に変換されてもよく、1 0, 0 0 0 個の疎要素は上に記載されるような第 M 番目の行 2 0 6 によってアクセス可能である。疎要素アクセスユニット  $X_{M,1}$  は、データベーステーブル 1, 0 0 0 番の疎要素 1 番 ~ 2, 0 0 0 番にアクセスするよう構成されてもよく、疎要素アクセスユニット  $X_{M,N}$  は、データベーステーブル 1, 0 0 0 番の疎要素 9, 0 0 0 番 ~ 1 0, 0 0 0 番にアクセスするよう構成されてもよい。

20

30

## 【 0 0 2 2 】

図 2 B は、疎密変換ユニット 2 0 0 が、疎要素アクセスユニットの二次元のメッシュネットワークを用いて、どのように疎要素を要求し得るかの例を示す。例として、処理ユニットは、疎密変換ユニット 2 0 0 に対して、データベーステーブル 1 番の疎要素 1 番 ~ 5 0 番、データベーステーブル 2 番の疎要素 1 0 0 番 ~ 2 0 0 番、およびデータベーステーブル 1, 0 0 0 番の疎要素 9, 0 5 0 番 ~ 9, 0 6 0 番を用いて生成される密な一次元のベクトルを要求する命令を実行してもよい。疎密変換ユニット 2 0 0 は、処理ユニットから要求を受けた後、疎密変換ユニット 2 0 0 は、疎要素アクセスユニット  $X_{1,1}$  に命令して、疎要素に対する要求をメッシュネットワークにおける他の疎要素アクセスユニットに同報通信させてもよい。疎要素アクセスユニット  $X_{1,1}$  は、疎要素アクセスユニット  $X_{1,2}$  に要求 2 2 2 を、および疎要素アクセスユニット  $X_{2,1}$  に要求 2 2 4 を同報通信してもよい。要求 2 2 2 を受けた後、疎要素アクセスユニット  $X_{1,2}$  は、要求 2 2 6 を疎要素アクセスユニット  $X_{1,3}$  に同報通信してもよい。いくつかの実現例では、疎要素アクセスユニットは、ルーティングスキームに基いて別の疎要素アクセスユニットに要求を同報通信するよう構成されてもよい。たとえば、疎要素アクセスユニット  $X_{1,2}$  は疎要素アクセスユニット  $X_{2,2}$  に要求を同報通信するよう構成されなくてもよく、なぜならば、疎要素

40

50

アクセスユニット $X_{2,2}$ は疎要素アクセスユニット $X_{2,1}$ から同報通信を受けるよう構成されるからである。ルーティングスキームは静的であってもよく、または動的に生成されてもよい。たとえば、ルーティングスキームはルックアップテーブルであってもよい。いくつかの実現例では、疎要素アクセスユニットは、要求224を要求224に基いて別の疎要素アクセスユニットに同報通信するよう構成されてもよい。たとえば、要求224は要求された疎要素の識別を含んでもよく（たとえばデータベーステーブル1番、疎要素1番～50番）、疎要素アクセスユニット $X_{1,2}$ は、要求224を疎要素アクセスユニット $X_{2,2}$ および/または疎要素アクセスユニット $X_{1,3}$ に同報通信すべきかどうかを、識別に基づいて判断してもよい。同報通信プロセスは、メッシュネットワークを介して伝搬し、疎要素アクセスユニット $X_{M,N}$ は疎要素アクセスユニット $X_{M,N-1}$ から要求230を受ける。

10

#### 【0023】

図2Cは、疎密変換ユニット200が、疎要素アクセスユニットの二次元のメッシュネットワークを用いて、要求される密行列をどのように生成し得るかの例を示す。いくつかの実現例では、或る疎要素アクセスユニットが同報通信された要求を受けた後、その疎要素アクセスユニットは、それは要求される疎要素のいずれかにアクセスするよう構成されるかどうかを判断するよう構成される。たとえば疎要素アクセスユニット $X_{1,1}$ は、それは、データベーステーブル1番の疎要素1番～50番にアクセスするよう構成されるが、データベーステーブル2番の疎要素100番～200番またはデータベーステーブル1,000番の疎要素9,050番～9,060番にアクセスするようには構成されない、と判断してもよい。それがデータベーステーブル1番の疎要素1番～50番にアクセスするよう構成されると判断することに応じて、疎要素アクセスユニット $X_{1,1}$ は、データベーステーブル1番の疎要素1番～50番を、これらの疎要素が格納されているデータ片からフェッチし、これらの疎要素に基いて密行列242を生成してもよい。

20

#### 【0024】

別の例として、疎要素アクセスユニット $X_{2,1}$ は、それが、データベーステーブル1番の疎要素1番～50番、データベーステーブル2番の疎要素100番～200番、またはデータベーステーブル1,000番の疎要素9,050番～9,060番のいずれにもアクセスするよう構成されない、と判断してもよい。それが要求される疎要素のいずれにもアクセスするよう構成されないと判断することに応じて、疎要素アクセスユニット $X_{2,1}$ はさらなるアクションを実行しなくてもよい。

30

#### 【0025】

別の例として、疎要素アクセスユニット $X_{1,2}$ は、それはデータベーステーブル2番の疎要素100番～200番にアクセスするよう構成されるが、データベーステーブル1番の疎要素1番～50番またはデータベーステーブル1,000番の疎要素9,050番～9,060番にアクセスするようには構成されない、と判断してもよい。それがデータベーステーブル2番の疎要素100番～200番にアクセスするよう構成されると判断することに応じて、疎要素アクセスユニット $X_{1,2}$ は、これらの疎要素が格納されているデータ片からこれらの疎要素をフェッチし、これらの疎要素に基いて密行列244を生成してもよい。いくつかの実現例では、ある疎要素アクセスユニットが、密行列を生成した後、その疎要素アクセスユニットは、同報通信された要求の送信側にその密行列を転送するよう構成されてもよい。ここでは、疎要素アクセスユニット $X_{1,2}$ は密行列244を疎要素アクセスユニット $X_{1,1}$ に転送する。

40

#### 【0026】

別の例として、疎要素アクセスユニット $X_{M,N}$ は、それはデータベーステーブル1,000番の疎要素9,050番～9,060番にアクセスするよう構成されるが、データベーステーブル1番の疎要素1番～50番またはデータベーステーブル2番の疎要素100番～200番にアクセスするようには構成されない、と判断してもよい。それがデータベーステーブル1,000番の疎要素9,050番～9,060番にアクセスするよう構成されると判断することに応じて、疎要素アクセスユニット $X_{M,N}$ は、これらの疎要素が格

50



納されているデータ片からこれらの疎要素をフェッチし、これらの疎要素に基いて密行列 2 4 6 を生成してもよい。いくつかの実現例では、ある疎要素アクセスユニットが、密行列を生成した後、その疎要素アクセスユニットは、同報通信された要求の送信側にその密行列を転送するよう構成されてもよい。ここでは、疎要素アクセスユニット  $X_{M, N}$  は密行列 2 4 6 を疎要素アクセスユニット  $X_{M, N-1}$  に転送する。次のサイクルで、疎要素アクセスユニット  $X_{M, N-1}$  は、密行列 2 4 6 を疎要素アクセスユニット  $X_{M, N-1}$  に転送するよう構成される。このプロセスは、疎要素アクセスユニット  $X_{2, 1}$  が密行列 2 4 6 を疎要素アクセスユニット  $X_{1, 1}$  に転送するまで継続する。

#### 【0027】

いくつかの実現例では、疎密変換ユニット 200 は、疎要素アクセスユニットによって生成された密行列を変換し、プロセッサユニットのための密行列を生成するように構成される。ここで、疎密変換ユニット 200 は、密行列 2 4 2、2 4 4 および 2 4 6 を、プロセッサユニットのための密行列に変換する。たとえば、密行列 2 4 2 は  $100 \times 10$  の次元を有してもよく、密行列 2 4 4 は  $20 \times 100$  の次元を有してもよく、密行列 2 4 6 は  $3 \times 3$  の次元を有してもよい。疎密変換ユニット 200 は、密行列 2 4 2、2 4 4 および 2 4 6 を、 $1 \times 3009$  の次元でベクトルに変換してもよい。有利なことに、密行列（たとえばデータベーステーブル）に従う行の区分化は、疎密変換ユニット 200 が、生成された密行列が列  $N$  から列 1 に伝搬した後に、要求された疎要素をすべて得ることを可能にする。列の区分化は、疎要素アクセスユニットのわずか 1 つを用いてあまりにも多数の疎要素にアクセスすることによって引起される帯域幅ボトルネックを低減する。

#### 【0028】

図 2 D は、疎密変換ユニット 200 が、疎要素アクセスユニットの二次元のメッシュネットワークを用いて、密行列に基いて疎要素をどのように更新し得るかの例を示す。例として、処理ユニットは、疎密変換ユニット 200 に対して、データベーステーブル 1 番の疎要素 1 番 ~ 50 番およびデータベーステーブル 1, 000 番の疎要素 9, 050 番 ~ 9, 060 番を用いて生成される密な一次元のベクトルを用いて、格納された疎要素を更新するよう要求する命令を実行してもよい。疎密変換ユニット 200 は、処理ユニットから要求を受けた後、疎密変換ユニット 200 は疎要素アクセスユニット  $X_{1, 1}$  に対して、疎要素更新要求をメッシュネットワークにおける他の疎要素アクセスユニットに同報通信するよう命令してもよく、疎要素更新要求は、処理ユニットによって与えられる密な一次元のベクトルを含んでもよい。いくつかの実現例では、疎要素アクセスユニット  $X_{1, 1}$  は、それが密な一次元のベクトルに含まれる疎要素にアクセスするよう割り当てられるかどうかを判断してもよい。それが密な一次元のベクトルに含まれる疎要素にアクセスするよう割り当てられると判断することに応じて、疎要素アクセスユニット  $X_{1, 1}$  は、データ片に格納される疎要素を更新してもよい。ここで、疎要素アクセスユニット  $X_{1, 1}$  は、それがデータベーステーブル 1 番の疎要素 1 番 ~ 50 番に割り当てられると判断し、疎要素アクセスユニット  $X_{1, 1}$  は、データ片におけるこれらの疎要素を更新するよう命令を実行する。

#### 【0029】

疎要素アクセスユニット  $X_{1, 1}$  は、疎要素アクセスユニット  $X_{1, 2}$  に疎要素更新要求 252 を、および疎要素アクセスユニット  $X_{2, 1}$  に疎要素更新要求 254 を同報通信してもよい。疎要素更新要求 252 を受けた後、疎要素アクセスユニット  $X_{1, 2}$  は、それは、密な一次元のベクトルに含まれる疎要素にアクセスするよう割り当てられない、と判断してもよい。疎要素アクセスユニット  $X_{1, 2}$  は、要求 256 を疎要素アクセスユニット  $X_{1, 3}$  に同報通信する。同報通信プロセスは、メッシュネットワークを介して伝搬し、疎要素アクセスユニット  $X_{M, N}$  は疎要素アクセスユニット  $X_{M, N-1}$  から要求 260 を受ける。ここで、疎要素アクセスユニット  $X_{M, N}$  は、それがデータベーステーブル 1, 000 番の疎要素 9, 050 番 ~ 9, 060 番に割り当てられると判断し、疎要素アクセスユニット  $X_{M, N}$  はデータ片におけるこれらの疎要素を更新するよう命令を実行する。

#### 【0030】

図 3 A は例示の疎要素アクセスユニット 300 を示す。疎要素アクセスユニット 300

は、疎要素アクセスユニット  $X_{1,1} \sim X_{M,N}$  の任意の 1 つであってもよい。一般に、疎要素アクセスユニット 300 は、ノードネットワーク 320 から、1 つ以上のデータ片に格納された疎要素をフェッチし、フェッチされた疎要素を密行列に変換する要求 342 を受けるよう構成される。いくつかの実現例では、処理ユニット 316 が、ノードネットワーク 320 における疎要素アクセスユニットに対して、疎要素を用いて生成された密行列を求める要求を送信する。疎要素アクセスユニットは、疎要素アクセスユニット 300 に要求 342 を同報通信してもよい。同報通信された要求 342 のルーティングは図 2 B における記載に類似してもよい。疎要素アクセスユニット 300 は、要求識別ユニット 302、データフェッチユニット 304、疎低減ユニット 306、連結ユニット 308、圧縮 / 伸長ユニット 310、および分割ユニット 312 を含む。ノードネットワーク 320 は二次元のメッシュネットワークであってもよい。処理ユニット 316 は処理ユニット 102 と類似してもよい。

10

#### 【0031】

一般に、要求識別ユニット 302 は、1 つ以上のデータ片 330 に格納された疎要素をフェッチするよう要求 342 を受け、疎要素アクセスユニット 300 は要求 342 によって示された疎要素にアクセスするよう割当てられるかどうかを判断するよう構成される。いくつかの実現例では、要求識別ユニット 302 は、疎要素アクセスユニット 300 が要求 342 によって示される疎要素にアクセスするよう割当てられるかどうかを、ルックアップテーブルを用いることによって判断してもよい。特定の要求された疎要素の識別（たとえば、データベーステーブル 1 番の 1 番）がルックアップテーブルに含まれている場合には、要求識別ユニット 302 は、特定の要求された疎要素をフェッチするよう、信号 344 をデータフェッチユニット 304 に送信してもよい。特定の要求された疎要素の識別（たとえば、データベーステーブル 1 番の 1 番）がルックアップテーブルに含まれない場合には、要求識別ユニット 302 は受取った要求を破棄してもよい。いくつかの実現例では、要求識別ユニット 302 は、受取った要求をノードネットワーク 320 上における別の疎要素アクセスユニットに同報通信するよう構成されてもよい。

20

#### 【0032】

データフェッチユニット 304 は、信号 344 を受信することに応じて、データ片 330 から 1 つ以上の要求された疎要素をフェッチするよう構成される。いくつかの実現例では、データフェッチユニット 304 は 1 つ以上のプロセッサ 322 a ~ 322 k を含み、 $k$  は整数である。プロセッサ 322 a ~ 322 k は、ベクトル処理ユニット (VPU)、アレイ処理ユニットまたは任意の好適な処理ユニットであってもよい。いくつかの実現例では、プロセッサ 322 a ~ 322 k は、データ片 330 近くに配置されて、プロセッサ 322 a ~ 322 k とデータ片 330 との間のレイテンシを低減するようにする。疎要素アクセスユニット 300 がフェッチするよう割当てられる、要求された疎要素の数に基いて、データフェッチユニット 304 は、プロセッサ 322 a ~ 322 k 間に分散されるべき 1 つ以上の要求を発生させるよう構成されてもよい。いくつかの実現例では、プロセッサ 322 a ~ 322 k の各々は、疎要素の識別に基いて特定の疎要素に割当てられてもよく、データフェッチユニット 304 は、プロセッサ 322 a ~ 322 k に対する 1 つ以上の要求を疎要素の識別に基いて発生させるよう構成されてもよい。いくつかの実現例では、データフェッチユニット 304 はルックアップテーブルを用いることによってプロセッサ割当てを判断してもよい。いくつかの実現例では、データフェッチユニット 304 は、プロセッサ 322 a ~ 322 k のために複数のパッチを生成してもよく、各パッチは要求された疎要素の部分集合に対する要求である。プロセッサ 322 a ~ 322 k は、割当てられた疎要素をデータ片 330 から独立してフェッチし、フェッチされた疎要素を疎低減ユニット 306 に転送するよう構成される。

30

40

#### 【0033】

疎低減ユニット 306 はフェッチされた疎要素 346 の次元を低減するように構成される。たとえば、プロセッサ 322 a ~ 322 k の各々は、 $100 \times 1$  の次元を有する疎要素を生成してもよい。疎低減ユニット 306 は、 $100 \times k$  の次元を有する、フェッチさ

50

れた疎要素 3 4 6 を受け、フェッチされた疎要素 3 4 6 の次元を論理演算、算術演算または両方の組合せによって 1 0 0 × 1 に低減することによって疎低減要素 3 4 8 を生成してもよい。疎低減ユニット 3 0 6 は疎低減要素 3 4 8 を連結ユニット 3 0 8 に出力するように構成される。

#### 【 0 0 3 4 】

連結ユニット 3 0 8 は、疎低減要素 3 4 8 を再配列および連結して、連結された要素 3 5 0 を生成するように構成される。たとえば、疎要素アクセスユニット X<sub>1</sub>,<sub>1</sub> は、データベーステーブル 1 番の疎要素 1 番 ~ 2 0 0 番にアクセスするよう構成されてもよい。プロセッサ 3 2 2 a は、フェッチされた疎要素 1 0 番を、フェッチされた疎要素 5 番を返すように構成されるプロセッサ 3 2 2 b よりも早く、疎低減ユニット 3 0 6 に返すかもしれない。連結ユニット 3 0 8 は、その後受取られる疎要素 5 番を、より早く受取られた疎要素 1 0 番の前に順序づけられるように再配列し、疎要素 1 番 ~ 2 0 0 番を連結された要素 3 5 0 として連結するよう構成される。

10

#### 【 0 0 3 5 】

圧縮 / 伸長ユニット 3 1 0 は、連結された要素 3 5 0 を圧縮して、ノードネットワーク 3 2 0 のための密行列 3 5 2 を生成するよう構成される。たとえば、圧縮 / 伸長ユニット 3 1 0 は、連結された要素 3 5 0 における零値を圧縮して、ノードネットワーク 3 2 0 の帯域幅を改善するよう構成されてもよい。いくつかの実現例では、圧縮 / 伸長ユニット 3 1 0 は、受取られた密行列を伸長してもよい。たとえば、疎要素アクセスユニット 3 0 0 は、ノードネットワーク 3 2 0 を介して近隣の疎要素アクセスユニットから密行列を受け

20

#### 【 0 0 3 6 】

図 3 B は、疎要素アクセスユニット 3 0 0 がノードネットワーク 3 2 0 から受取られる密行列に基いて疎要素をどのように更新し得るかの例を示す。例として、処理ユニットは、疎密変換ユニットに対して、データベーステーブル 1 番の疎要素 1 番 ~ 5 0 番およびデータベーステーブル 1, 0 0 0 番の疎要素 9, 0 5 0 番 ~ 9, 0 6 0 番を用いて生成される密な一次元のベクトルを用いて、格納された疎要素を更新するよう要求する命令を実行してもよい。疎密変換ユニットは、処理ユニットから要求を受けた後、疎密変換ユニットは要求 3 6 2 を送信して、疎要素アクセスユニット 3 0 0 に対して、それが密な一次元のベクトルに含まれる疎要素にアクセスするよう割り当てられるかどうかを判断するよう命令してもよい。要求識別ユニット 3 0 2 は、疎要素アクセスユニット 3 0 0 が密な一次元のベクトルに含まれる疎要素にアクセスするよう割り当てられるかどうかを判断するよう構成される。疎要素アクセスユニット 3 0 0 が密な一次元のベクトルに含まれる疎要素にアクセスするよう割り当てられると判断することに応じて、要求識別ユニット 3 0 2 は、データ片において格納された疎要素を更新するよう、指示 3 6 4 を分割ユニット 3 1 2 に送信してもよい。

30

#### 【 0 0 3 7 】

分割ユニット 3 1 2 は、受取られた密行列を、データ片 3 3 0 においてデータフェッチユニット 3 0 4 によって更新することができる疎要素に変換するように構成される。たとえば、分割ユニット 3 1 2 は、密な一次元のベクトルを複数の疎要素に変換し、データフェッチユニット 3 0 4 に対して、疎要素アクセスユニット 3 0 0 がフェッチするよう割り当てられるデータ片 3 3 0 において格納された疎要素を更新するよう命令するよう構成されてもよい。

40

#### 【 0 0 3 8 】

図 4 は、密行列を生成するためのプロセス 4 0 0 の例を示すフローチャートである。プロセス 4 0 0 は、疎密変換ユニット 1 0 4 または疎密変換ユニット 2 0 0 などのようなシステムによって実行されてもよい。システムは、疎要素アクセスユニットの第 1 の群および疎要素アクセスユニットの第 2 の群を含んでもよい。たとえば、図 2 A を参照して、疎

50

疎変換ユニット 200 は、 $M \times N$  個の疎要素アクセスユニット  $X_{1,1} \sim X_{M,N}$  を含み、それらは物理的にまたは論理的に  $M$  個の行および  $N$  個の列に配列される。疎要素アクセスユニット  $X_{1,1} \sim X_{M,N}$  の各行は、特定の密行列から変換される疎要素にアクセスするよう区分されてもよい。いくつかの実現例では、疎要素アクセスユニットの第 1 の群は第 1 の疎要素アクセスユニットおよび第 2 の疎要素アクセスユニットを含んでもよい。たとえば、疎変換ユニット 200 の第 1 番目の行は疎要素アクセスユニット  $X_{1,1}$  および、 $X_{1,2}$  を含んでもよい。いくつかの実現例では、疎要素アクセスユニットの第 1 の群および疎要素アクセスユニットの第 2 の群は二次元のメッシュ構成で配列されてもよい。いくつかの実現例では、疎要素アクセスユニットの第 1 の群および疎要素アクセスユニットの第 2 の群は二次元の円環面構成で配列されてもよい。

10

#### 【0039】

システムは、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を含む疎要素に基いた出力行列に対する要求を受ける (402)。たとえば、図 2B を参照して、処理ユニットは、疎変換ユニット 200 に対して、データベーステーブル 1 番の疎要素 1 番 ~ 50 番、データベーステーブル 2 番の疎要素 100 番 ~ 200 番およびデータベーステーブル 1,000 番の疎要素 9,050 番 ~ 9,060 番を用いて生成される密な一次元のベクトルを要求する命令を実行してもよい。

#### 【0040】

いくつかの実現例では、第 1 の疎要素アクセスユニットは、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を含む複数の疎要素に対する要求を受取ってもよい。第 1 の疎要素アクセスユニットは、第 2 の疎要素アクセスユニットに要求を送信してもよい。たとえば、図 2B を参照して、疎変換ユニット 200 は、処理ユニットから要求を受けた後、疎変換ユニット 200 は、疎要素アクセスユニット  $X_{1,1}$  に命令して、疎要素に対する要求をメッシュネットワークにおける他の疎要素アクセスユニットに同報通信させてもよい。疎要素アクセスユニット  $X_{1,1}$  は、要求 222 を疎要素アクセスユニット  $X_{1,2}$  に同報通信してもよい。

20

#### 【0041】

システムは、疎要素アクセスユニットの第 1 の群によってフェッチされる第 1 の密行列と関連付けられる疎要素を得る (404)。いくつかの実現例では、第 1 の疎要素アクセスユニットは、複数の疎要素のうちの特定の疎要素のアイデンティティが、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合のうちの 1 つのアイデンティティと一致する、と判断してもよい。たとえば、図 2C を参照して、疎要素アクセスユニット  $X_{1,1}$  は、データベーステーブル 1 番の疎要素 1 番 ~ 200 番にアクセスするよう構成されてもよい。疎要素アクセスユニット  $X_{1,1}$  は、それはデータベーステーブル 1 番の疎要素 1 番 ~ 50 番にアクセスするよう構成されるが、データベーステーブル 2 番の疎要素 100 番 ~ 200 番またはデータベーステーブル 1,000 番の疎要素 9,050 番 ~ 9,060 番にアクセスするようには構成されない、と判断してもよい。複数の疎要素のうちの特定の疎要素のアイデンティティが、第 1 の密行列と関連付けられる疎要素の第 1 の部分集合のうちの 1 つのアイデンティティと一致する、と判断することに応じて、第 1 の疎要素アクセスユニットは、特定の疎要素を含む第 1 の密行列と関連付けられる疎要素の第 1 の部分集合をフェッチしてもよい。たとえば、それがデータベーステーブル 1 番の疎要素 1 番 ~ 50 番にアクセスするよう構成されると判断することに応じて、疎要素アクセスユニット  $X_{1,1}$  は、データベーステーブル 1 番の疎要素 1 番 ~ 50 番を、これらの疎要素が格納されているデータ片からフェッチしてもよい。

30

40

#### 【0042】

第 2 の疎要素アクセスユニットは、第 1 の密行列と関連付けられる疎要素の第 2 の異なる部分集合をフェッチしてもよい。たとえば、図 2C を参照して、疎要素アクセスユニット  $X_{1,2}$  は、データベーステーブル 2 番の疎要素 51 番 ~ 200 番にアクセスするよう構成されてもよい。それがデータベーステーブル 2 番の疎要素 100 番 ~ 200 番にアクセスするよう構成されると判断することに応じて、疎要素アクセスユニット  $X_{1,2}$  は、

50

これらの疎要素が格納されているデータ片からこれらの疎要素をフェッチしてもよい。

#### 【 0 0 4 3 】

システムは、疎要素アクセスユニットの第 2 の群によってフェッチされる第 2 の密行列と関連付けられる疎要素を得る ( 4 0 6 )。たとえば、図 2 C を参照して、第 2 の群疎要素アクセスユニットは、 $M \times N$  個の疎要素アクセスユニットの  $M$  番目の行であってもよく、疎要素アクセスユニット  $X_{M, N}$  は、データベーステーブル 1, 0 0 0 番の疎要素 9, 0 0 0 番 ~ 1 0, 0 0 0 番にアクセスするよう構成されてもよい。それがデータベーステーブル 1, 0 0 0 番の疎要素 9, 0 5 0 番 ~ 9, 0 6 0 番にアクセスするよう構成されると判断することに応じて、疎要素アクセスユニット  $X_{M, N}$  は、これらの疎要素が格納されているデータ片からこれらの疎要素をフェッチし、これらの疎要素に基づいて密行列 2 4 6 を生成してもよい。

10

#### 【 0 0 4 4 】

いくつかの実現例では、第 1 の疎要素アクセスユニットは、第 1 のデータ片から第 1 の密行列と関連付けられる疎要素の第 1 の部分集合をフェッチしてもよく、第 2 の疎要素アクセスユニットは、第 2 の異なるデータ片から第 1 の密行列と関連付けられる疎要素の第 2 の異なる部分集合をフェッチしてもよい。たとえば、図 1 を参照して、第 1 の疎要素アクセスユニットは、データ片 1 0 6 a から第 1 の密行列と関連付けられる疎要素の第 1 の部分集合をフェッチしてもよく、第 2 の疎要素アクセスユニットは、データ片 1 0 6 b から第 1 の密行列と関連付けられる疎要素の第 2 の異なる部分集合をフェッチしてもよい。

#### 【 0 0 4 5 】

20

システムは、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を変換して、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素を含む出力密行列を生成する ( 4 0 8 )。たとえば、図 2 C を参照して、疎密変換ユニット 2 0 0 は、密行列 2 4 2、2 4 4 および 2 4 6 を、プロセッサユニットのための密行列に変換してもよい。

#### 【 0 0 4 6 】

いくつかの実現例では、第 1 の密行列と関連付けられる疎要素および第 2 の密行列と関連付けられる疎要素は多次元の行列であってもよく、出力密行列はベクトルであってもよい。たとえば、密行列 2 4 2 は  $1 0 0 \times 1 0$  の次元を有してもよく、密行列 2 4 4 は  $2 0 \times 1 0 0$  の次元を有してもよく、密行列 2 4 6 は  $3 \times 3$  の次元を有してもよい。疎密変換ユニット 2 0 0 は、密行列 2 4 2、2 4 4 および 2 4 6 を、 $1 \times 3 0 0 9$  の次元でベクトルに変換してもよい。

30

#### 【 0 0 4 7 】

図 5 は、密行列を生成するためのプロセス 5 0 0 の例を示すフローチャートである。プロセス 5 0 0 は、疎密変換ユニット 1 0 4 または疎要素アクセスユニット 3 0 0 などのようなシステムによって実行されてもよい。

#### 【 0 0 4 8 】

システムは特定の疎要素の部分集合にアクセスすることに対する指示を受ける ( 5 0 2 )。たとえば図 3 A を参照して、データフェッチユニット 3 0 4 は、データ片 3 3 0 から 1 つ以上の要求された疎要素をフェッチするための信号 3 4 4 を受信するよう構成されてもよい。いくつかの実現例では、1 つ以上のデータ片において格納される特定の疎要素に対する要求が、ノードネットワーク上で受けられてもよい。たとえば図 3 A を参照して、要求識別ユニット 3 0 2 は、データ片 3 3 0 において格納された疎要素をフェッチするようノードネットワーク 3 2 0 上で要求 3 4 2 を受けるよう構成されてもよい。システムは、データフェッチユニットは特定の疎要素の部分集合を扱うよう割当てられる、と判断してもよい。たとえば、要求識別ユニット 3 0 2 は、疎要素アクセスユニット 3 0 0 は要求 3 4 2 によって示される疎要素にアクセスするよう割当てられるかどうかを判断するよう構成されてもよい。データフェッチユニットは特定の疎要素の部分集合を扱うよう割当てられる、と判断することに応じて、その指示は特定の疎要素の部分集合にアクセスすることに対して生成されてもよい。たとえば、特定の要求された疎要素の識別 (たとえば、デ

40

50

ータベーステーブル 1 番の 1 番) がルックアップテーブルに含まれている場合には、要求識別ユニット 3 0 2 は、特定の要求された疎要素をフェッチするよう、信号 3 4 4 をデータフェッチユニット 3 0 4 に送信してもよい。

【 0 0 4 9 】

システムは、特定の疎要素の部分集合の識別に基いて、特定の疎要素の部分集合をフェッチするためのプロセッサ指定を判断する ( 5 0 4 )。たとえば図 3 A を参照して、データフェッチユニット 3 0 4 は 1 つ以上のプロセッサ 3 2 2 a ~ 3 2 2 k を含む。プロセッサ 3 2 2 a ~ 3 2 2 k の各々は、疎要素の識別に基いて特定の疎要素に割当てられてもよく、データフェッチユニット 3 0 4 は、プロセッサ 3 2 2 a ~ 3 2 2 k に対する 1 つ以上の要求を疎要素の識別に基いて発生させるよう構成されてもよい。いくつかの実現例では、システムは、システムが特定の疎要素の部分集合を扱うよう割当てられると判断してもよく、システムはルックアップテーブルに基いて特定の疎要素の部分集合を扱うよう割当てられると判断することを含む。たとえば、データフェッチユニット 3 0 4 はルックアップテーブルを用いることによってプロセッサ割当を判断してもよい。

10

【 0 0 5 0 】

システムは、指定に基いて、および複数個のプロセッサのうちの第 1 のプロセッサによって、特定の疎要素の部分集合の第 1 の疎要素をフェッチする ( 5 0 6 )。たとえば図 3 A を参照して、データフェッチユニット 3 0 4 は、信号 3 4 4 に含まれる疎要素をフェッチするようプロセッサ 3 2 2 a に命令してもよい。

【 0 0 5 1 】

20

システムは、指定に基いて、および複数個のプロセッサのうちの第 2 のプロセッサによって、特定の疎要素の部分集合の第 2 の疎要素をフェッチする ( 5 0 8 )。たとえば図 3 A を参照して、データフェッチユニット 3 0 4 は、信号 3 4 4 に含まれる異なる疎要素をフェッチするようプロセッサ 3 2 2 b に命令してもよい。

【 0 0 5 2 】

いくつかの実現例では、第 1 のプロセッサから第 1 の疎要素を含む第 1 の行列を受取ってもよく、第 1 の行列は第 1 の次元を有してもよい。システムは、第 1 の疎要素を含む第 2 の行列を生成してもよく、第 2 の行列は、第 1 の次元よりも小さい第 2 の次元を有する。たとえば、疎低減ユニット 3 0 6 はフェッチされた疎要素 3 4 6 の次元を低減するように構成されてもよい。プロセッサ 3 2 2 a ~ 3 2 2 k の各々は、 $100 \times 1$  の次元を有する疎要素を生成してもよい。疎低減ユニット 3 0 6 は、 $100 \times k$  の次元を有する、フェッチされた疎要素 3 4 6 を受け、フェッチされた疎要素 3 4 6 の次元を論理演算、算術演算または両方の組合せによって  $100 \times 1$  に低減することによって疎低減要素 3 4 8 を生成してもよい。システムは出力密行列を生成してもよく、出力密行列は第 2 の行列に基いて生成されてもよい。たとえば、連結ユニット 3 0 8 は、疎低減要素 3 4 8 を再配列および連結して、連結された要素 3 5 0 を生成するように構成されてもよい。

30

【 0 0 5 3 】

いくつかの実現例では、第 1 の疎要素は第 1 の時間の点において受取られてもよく、第 2 の疎要素は第 2 の異なる時間の点において受取られてもよい。システムは、出力密行列のために第 1 の疎要素および第 2 の疎要素の順序を判断してもよい。たとえば図 3 A を参照して、プロセッサ 3 2 2 a は、フェッチされた疎要素 1 0 番を、フェッチされた疎要素 5 番を返すように構成されるプロセッサ 3 2 2 b よりも早く、疎低減ユニット 3 0 6 に返すかもしれない。連結ユニット 3 0 8 は、その後受取られる疎要素 5 番を、より早く受取られた疎要素 1 0 番の前に順序づけられるように再配列し、疎要素 1 番 ~ 2 0 0 番を連結された要素 3 5 0 として連結するよう構成される。

40

【 0 0 5 4 】

システムは、少なくとも第 1 の疎要素および第 2 の疎要素に適用される変換に基いて出力密行列を生成する ( 5 1 0 )。いくつかの実現例では、システムは、出力密行列を圧縮して、圧縮された出力密行列を生成してもよい。システムは、圧縮された出力密行列をノードネットワークに与えてもよい。たとえば、圧縮 / 伸長ユニット 3 1 0 は、連結された

50

要素 3 5 0 を圧縮して、ノードネットワーク 3 2 0 のための密行列 3 5 2 を生成するよう構成されてもよい。

【 0 0 5 5 】

いくつかの実現例では、システムは、ノードネットワーク上で送信される密行列を表す第 1 の密行列を受取ってもよく、第 1 の密行列、第 1 の疎要素および第 2 の疎要素に基いて出力密行列を生成してもよい。たとえば、疎要素アクセスユニット 3 0 0 は、ノードネットワーク 3 2 0 を介して近隣の疎要素アクセスユニットから密行列を受けてもよい。疎要素アクセスユニット 3 0 0 は受取られた密行列を伸長してもよく、伸長された密行列を連結された要素 3 5 0 と連結して、更新された連結された要素を形成してもよく、それらは圧縮され、次いでノードネットワーク 3 2 0 に出力されることができる。

10

【 0 0 5 6 】

いくつかの実現例では、特定の疎要素のうちの 1 つ以上の疎要素は、多次元の行列であり、出力密行列はベクトルである。主題の実施形態および本明細書に記載される機能的動作は、デジタル電子回路において、有形に実施されたコンピュータソフトウェアもしくはファームウェアにおいて、本明細書において開示された構造およびそれらの構造等価物を含むコンピュータソフトウェアにおいて、または、それらの 1 つ以上の組合せにおいて実現され得る。本明細書に記載される主題の実施形態は、1 つ以上のコンピュータプログラムとして、すなわち、データ処理装置による実行のために、または、データ処理装置の動作を制御するために有形の非一時的なプログラム担体上でエンコードされたコンピュータプログラム命令の 1 つ以上のモジュールとして実現され得る。代替的に、または加えて、プログラム命令は、データ処理装置による実行に対して好適な受信側装置への送信のために情報をエンコードするよう生成される、たとえばマシンにより生成された電気信号、光信号、または電磁気信号などの、人為的に生成された伝播される信号上でエンコードすることができる。コンピュータ記憶媒体は、コンピュータ読取可能記憶装置、コンピュータ読取可能記憶基板、ランダムもしくはシリアルアクセスメモリデバイス、または、それらの 1 つ以上の組合せであり得る。

20

【 0 0 5 7 】

「データ処理装置」という用語は、例としてプログラマブルプロセッサ、コンピュータ、または複数のプロセッサもしくはコンピュータを含む、データを処理するためのすべての種類の装置、デバイスおよびマシンを包含する。当該装置は、たとえば F P G A ( フィールドプログラマブルゲートアレイ ) または A S I C ( 特定用途向け集積回路 ) といった特定目的論理回路を含み得る。当該装置は、ハードウェアに加えて、たとえばプロセッサファームウェア、プロトコルスタック、データベース管理システム、オペレーティングシステム、または、それらの 1 つ以上の組合せを構成するコードといった、当該コンピュータプログラムについて実行環境を作成するコードをさらに含み得る。

30

【 0 0 5 8 】

( プログラム、ソフトウェア、ソフトウェアアプリケーション、モジュール、ソフトウェアモジュール、スクリプトまたはコードとも称され、または記載され得る ) コンピュータプログラムは、コンパイル型もしくはインタープリタ型言語、または宣言型もしくは手続き型言語を含む任意の形態のプログラミング言語で記述され得、スタンドアロンプログラムとして、または、モジュール、コンポーネント、サブルーチン、もしくは、コンピューティング環境で使用するのに好適な他のユニットとして任意の形態で展開され得る。コンピュータプログラムは、ファイルシステムにおけるファイルに対応し得るが、対応する必要があるわけではない。プログラムは、当該プログラムに専用である単一のファイルにおいて、または、複数の連携ファイル ( coordinated files ) ( たとえばコードの 1 つ以上のモジュール、サブプログラムまたは部分を格納するファイル ) において、他のプログラムまたはデータ ( たとえばマークアップ言語ドキュメントに格納される 1 つ以上のスクリプト ) を保持するファイルの一部に格納され得る。コンピュータプログラムは、1 つの場所に位置するかもしくは複数の場所にわたって分散され通信ネットワークによって相互接続される 1 つのコンピュータまたは複数のコンピュータ上で実行されるように展開され

40

50

得る。

【 0 0 5 9 】

本明細書に記載されるプロセスおよび論理フローは、入力データ上で動作し出力を生成することにより機能を実行するよう1つ以上のプログラマブルプロセッサが1つ以上のコンピュータプログラムを実行することによって実行され得る。プロセスおよび論理フローは、たとえばFPGA（フィールドプログラマブルゲートアレイ）、ASIC（特定用途向け集積回路）といった特殊目的論理回路として、またはGPGPU（汎用グラフィック処理装置）として実現され得る。

【 0 0 6 0 】

コンピュータプログラムの実行に好適であるプロセッサは、例として、汎用マイクロプロセッサもしくは特殊目的マイクロプロセッサもしくはその両方または任意の種類の中央処理ユニットに基づき得る。一般に、中央処理ユニットは、リードオンリメモリもしくはランダムアクセスメモリまたはその両方から命令およびデータを受取る。コンピュータの必須の要素は、命令を実行するための中央処理ユニットと、命令およびデータを格納するための1つ以上のメモリデバイスとである。一般に、コンピュータはさらに、たとえば磁気ディスク、光磁気ディスクまたは光ディスクといった、データを格納するための1つ以上の大容量記憶装置を含むか、当該1つ以上の大容量記憶装置からデータを受取るかもしくは当該1つ以上の大容量記憶装置にデータを転送するよう動作可能に結合されるか、またはその両方を行う。しかしながら、コンピュータはそのような装置を有する必要はない。さらに、コンピュータはたとえば、携帯電話、携帯情報端末（PDA）、モバイルオーディオまたはビデオプレーヤ、ゲームコンソール、全地球測位システム（GPS）受信機、またはポータブル記憶装置（たとえばユニバーサルシリアルバス（USB）フラッシュドライブ）といった別のデバイスに埋め込まれ得る。

【 0 0 6 1 】

コンピュータプログラム命令およびデータを格納するのに好適であるコンピュータ読取可能媒体は、例として、たとえばEPROM、EEPROMおよびフラッシュメモリ素子といった半導体メモリデバイスと、たとえば内部ハードディスクまたはリムーバブルディスクといった磁気ディスクと、光磁気ディスクと、CD-ROMおよびDVD-ROMディスクとを含むすべての形態の不揮発性メモリ、媒体およびメモリデバイスを含む。プロセッサおよびメモリは、特殊目的論理回路によって補足され得るか、または特殊目的論理回路に組み込まれ得る。

【 0 0 6 2 】

ユーザとのインタラクションを提供するために、本明細書に記載される主題の実施形態は、たとえばCRT（陰極線管）またはLCD（液晶ディスプレイ）モニタといったユーザに対して情報を表示するための表示デバイスと、たとえばマウス、トラックボールといったユーザがコンピュータに入力を提供可能であるキーボードおよびポインティングデバイスとを有するコンピュータ上で実現され得る。他の種類のデバイスが同様に、ユーザとのインタラクションを提供するために使用され得；たとえば、ユーザに提供されるフィードバックは、たとえば視覚フィードバック、聴覚フィードバックまたは触覚フィードバックといった任意の形態の感覚フィードバックであり得；ユーザからの入力は、音響入力、音声入力、または触覚入力を含む任意の形態で受取られ得る。加えて、コンピュータは、ユーザが使用するデバイスにドキュメントを送信しユーザが使用するデバイスからドキュメントを受信することによって、たとえば、ウェブブラウザから受信された要求にตอบสนองしてユーザのクライアントデバイス上のウェブブラウザにウェブページを送信することによって、ユーザと対話し得る。

【 0 0 6 3 】

本明細書に記載される主題の実施形態は、たとえばデータサーバとしてバックエンドコンポーネントを含む計算システムにおいて実現され得るか、たとえばアプリケーションサーバといったミドルウェアコンポーネントを含む計算システムにおいて実現され得るか、たとえば本明細書に記載される主題の実現例とユーザが対話することが可能であるグラフ

10

20

30

40

50



ィカルユーザーインターフェイスもしくはウェブブラウザを有するクライアントコンピュータといったフロントエンドコンポーネントを含む計算システムにおいて実現され得るか、または1つ以上のそのようなバックエンドコンポーネント、ミドルウェアコンポーネントもしくはフロントエンドコンポーネントの任意の組合せの計算システムにおいて実現され得る。システムのコンポーネントは、たとえば通信ネットワークといったデジタルデータ通信の任意の形態または媒体によって相互接続され得る。通信ネットワークの例は、ローカルエリアネットワーク(「LAN」)およびワイドエリアネットワーク(「WAN」)、たとえばインターネットを含む。

#### 【0064】

計算システムはクライアントおよびサーバを含むことができる。クライアントとサーバとは一般に互いから遠隔にあり、典型的には通信ネットワークを通じて対話する。クライアントとサーバとの関係は、それぞれのコンピュータ上で実行されるとともに互いに対してクライアント-サーバ関係を有するコンピュータプログラムによって発生する。

#### 【0065】

本明細書は多くの特定の實現例の詳細を含んでいるが、これらは如何なる発明の範囲または請求され得るものの範囲に対する限定としても解釈されるべきではなく、特定の発明の特定の實施形態に特有の特徴であり得る記載として解釈されるべきである。別個の實施形態の文脈で本明細書において記載されるある特徴は、単一の實施形態において組合せでも實現され得る。反対に、単一の實施形態の文脈において記載されるさまざまな特徴は、複数の實施形態において別々に、または任意の好適な部分的組合せでも實現され得る。さらに、特徴は、ある組合せにおいて作用すると上で記載されているとともにそのように最初は請求されている場合があるが、請求される組合せのうちの1つ以上の特徴はいくつかの場合には当該組合せから削除され得、請求される組合せは、部分的組合せまたは部分的組合せの変形例に関し得る。

#### 【0066】

同様に、動作が図においては特定の順に示されているが、そのような動作は、望ましい結果を達成するために、示された当該特定の順もしくは連続した順で実行される必要があると理解されるべきではなく、または、すべての示された動作が実行される必要があると理解されるべきではない。ある状況においては、マルチタスキングおよび並列処理が有利であり得る。さらに、上に記載された實施形態におけるさまざまなシステムモジュールおよびコンポーネントの分離は、すべての實施形態においてそのような分離を必要とするとは理解されるべきでなく、記載されたプログラムコンポーネントおよびシステムは一般に、単一のソフトウェアプロダクトと一緒に統合され得るか、または、複数のソフトウェアプロダクトへとパッケージ化され得るということが理解されるべきである。

#### 【0067】

主題の特定の實施形態が記載された。他の實施形態は以下の請求の範囲内にある。たとえば、請求項において記載されるアクションは、異なる順で実行され得、それでも望ましい結果を達成し得る。一例として、添付の図において示されるプロセスは、望ましい結果を達成するために、示された特定の順または連続する順であることを必ずしも必要としない。ある實現例においては、マルチタスキングおよび並列処理が有利であり得る。

#### 【符号の説明】

#### 【0068】

100 計算システム、102 処理ユニット、104 疎密変換ユニット、106 a ~ 106 k データ片、200 疎密変換ユニット、300 疎要素アクセスユニット、302 要求識別ユニット、304 データフェッチユニット、306 疎低減ユニット、308 連結ユニット、310 圧縮/伸長ユニット、312 分割ユニット、330 データ片、320 ノードネットワーク、322 a ~ 322 k プロセッサ。

10

20

30

40

【図面】

【図 1】

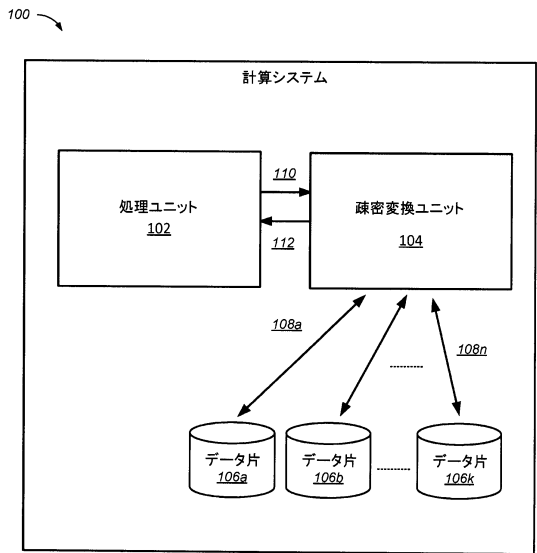


FIG. 1

【図 2 A】

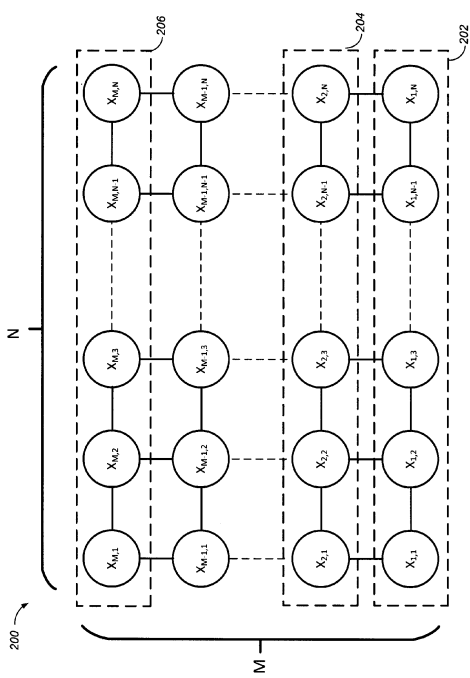


FIG. 2A

【図 2 B】

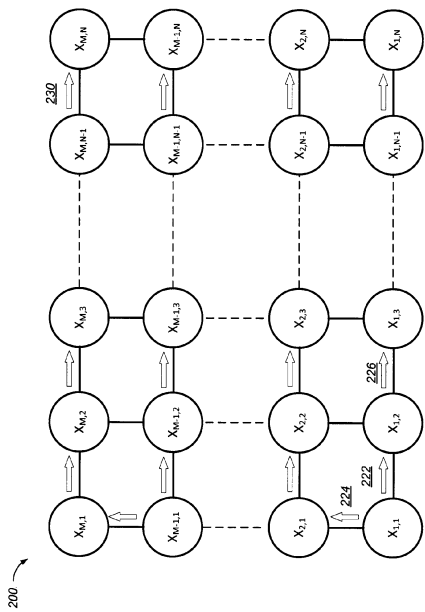


FIG. 2B

【図 2 C】

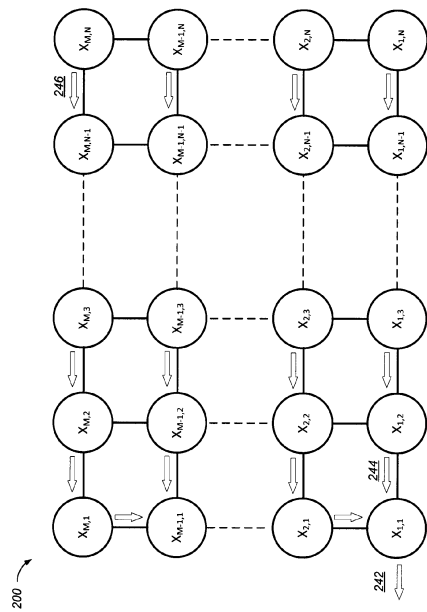


FIG. 2C

10

20

30

40

50

【図 2 D】

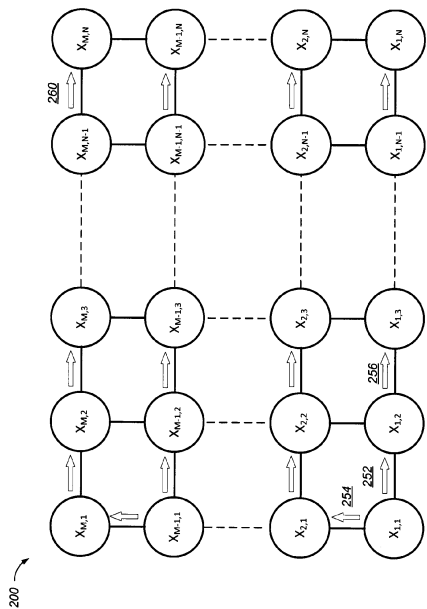


FIG. 2D

【図 3 A】

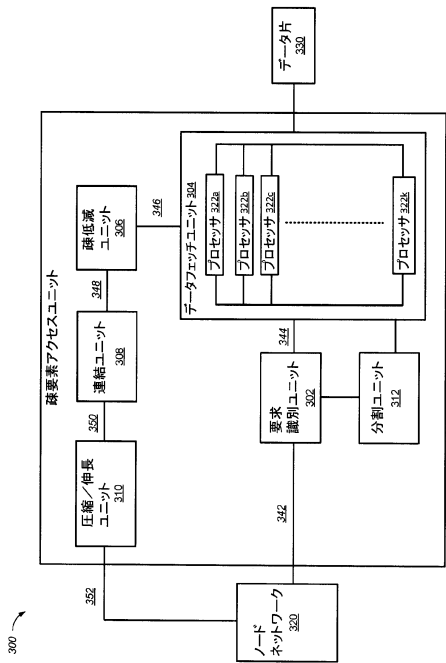


FIG. 3A

【図 3 B】

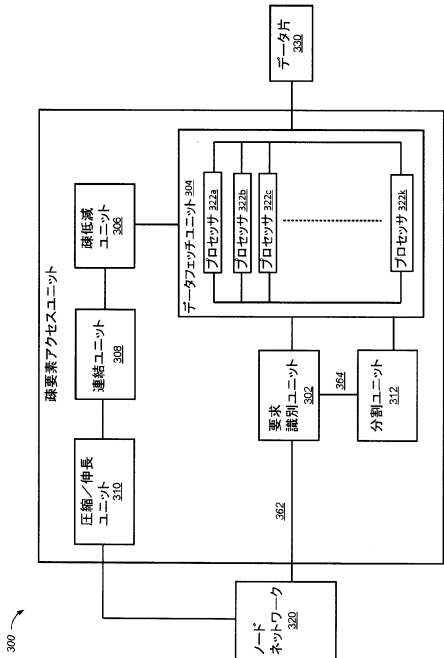


FIG. 3B

【図 4】

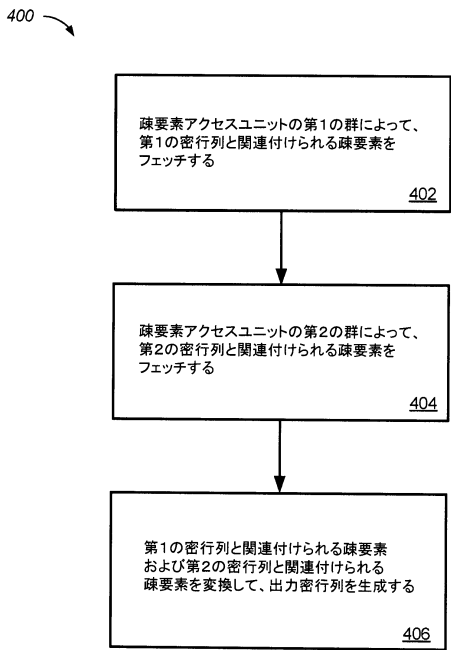


FIG. 4

10

20

30

40

50

【 図 5 】

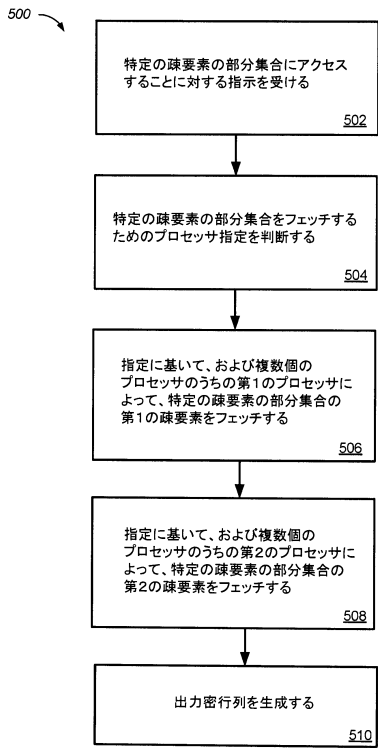


FIG. 5

## フロントページの続き

- (72)発明者    ラフル・ナガラジャン  
                 アメリカ合衆国、 9 4 0 4 3   カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、 1 6 0 0
- (72)発明者    ウ・ドン・ヒョク  
                 アメリカ合衆国、 9 4 0 4 3   カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、 1 6 0 0
- (72)発明者    クリストファー・ダニエル・リアリー  
                 アメリカ合衆国、 9 4 0 4 3   カリフォルニア州、マウンテン・ビュー、アンフィシアター・パークウェイ、 1 6 0 0
- 審査官    坂庭   剛史
- (56)参考文献    米国特許出願公開第 2 0 1 5 / 0 0 6 7 0 0 9 ( U S , A 1 )  
                 特開 2 0 1 4 - 1 9 9 5 4 5 ( J P , A )  
                 特開平 0 6 - 0 5 2 1 2 5 ( J P , A )
- (58)調査した分野 (Int.Cl. , D B 名)  
                 G 0 6 F    1 7 / 1 6