

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7386901号
(P7386901)

(45)発行日 令和5年11月27日(2023.11.27)

(24)登録日 令和5年11月16日(2023.11.16)

(51)国際特許分類 F I
G 0 6 F 17/16 (2006.01) G 0 6 F 17/16 B

請求項の数 23 (全31頁)

(21)出願番号	特願2021-566236(P2021-566236)	(73)特許権者	390009531
(86)(22)出願日	令和2年5月1日(2020.5.1)		インターナショナル・ビジネス・マシ
(65)公表番号	特表2022-531786(P2022-531786		ンズ・コーポレーション
	A)		INTERNATIONAL BUSI
(43)公表日	令和4年7月11日(2022.7.11)		NESS MACHINES CORPO
(86)国際出願番号	PCT/IB2020/054137		RATION
(87)国際公開番号	WO2020/229933		アメリカ合衆国10504 ニューヨー
(87)国際公開日	令和2年11月19日(2020.11.19)		ク州 アーモンク ニュー オーチャード
審査請求日	令和4年10月21日(2022.10.21)		ロード
(31)優先権主張番号	16/408,575		New Orchard Road, A
(32)優先日	令和1年5月10日(2019.5.10)		rmonk, New York 105
(33)優先権主張国・地域又は機関	米国(US)	(74)代理人	04, United States of
			America
			100112690
			弁理士 太佐 種一

最終頁に続く

(54)【発明の名称】 高性能なベクトル処理のためのアドレス生成

(57)【特許請求の範囲】

【請求項1】

バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んで
いるメモリ・ユニット内の前記バイナリ・データ・ベクトルにアクセスするための方法で
あって、前記方法が、

前記バイナリ・データ・ベクトルの開始アドレスおよび前記バイナリ・データ・ベクト
ルの要素の2の累乗のストライドsを受信することと、

前記メモリ・バンクの各々について1つのn個のオフセットを決定することであって、
前記オフセットの各々が、

複数のビットレベルのXOR関数を前記開始アドレスに適用してZベクトルを生成する
ことと、

マッピング・テーブルにアクセスするために前記Zベクトルを使用することと、

前記バイナリ・データ・ベクトルの2の累乗のストライドに従ってマッピング・テーブ
ルのアクセス結果をシフトすることによって決定される、前記決定することと、

前記Zベクトルの2進等価値に応じて前記n個のメモリ・バンク内の前記バイナリ・デ
ータ・ベクトルの部分のシーケンスを決定することと、

前記メモリ・ユニットの前記n個のメモリ・バンク内の前記バイナリ・データ・ベクト
ルに並列にアクセスすることを含む、方法。

【請求項2】

複数のビットレベルのXOR関数を前記開始アドレスに前記適用することが、

10

20

前記開始アドレスから $\log_2(n)$ 個の部分および $\log_2(n)$ ビット・サイズの Y の部分を選択することと、

前記 $\log_2(n)$ 個の部分の各々についてパリティ・ビットを決定し、 $\log_2(n)$ 次元のパリティ・ベクトル X を生成することと、

X と Y の間でビットレベルの XOR 演算を実行し、 Z ベクトルを生成することを含む、請求項 1 に記載の方法。

【請求項 3】

マッピング・テーブルにアクセスするために前記 Z ベクトルを使用することが、

前記マッピング・テーブル内のインデックスとして前記 Z ベクトルを使用することによって、 $n * n$ 個のエントリを含んでいる前記マッピング・テーブル内の行を選択することと、

10

前記マッピング・テーブルのアクセス結果を前記開始アドレスの部分と結合して、前記 n 個のオフセットを取得することとも含む、請求項 2 に記載の方法。

【請求項 4】

前記結合することが、

前記開始アドレスのビット・オフセット $\log_2(n)$ で開始する $\log_2(n)$ ビットを、前記決定された $\log_2(n)$ 個のマッピング・テーブル・エントリに置き換え、 n 個のオフセットを生成することを含み、各オフセットが、前記 n 個のバンクの各 1 つにおいてオフセットとして使用される、請求項 3 に記載の方法。

【請求項 5】

20

前記 n 個のメモリ・バンク内の前記バイナリ・データ・ベクトルの前記部分のシーケンスを前記決定することが、

$\log_2(n)$ 個の連続するマルチプレクサ段を制御することを含み、前記マルチプレクサ段の各々が、前記バイナリ・データ・ベクトルの全幅をカバーするように、複数の 2 入力ビット・マルチプレクサを含む、請求項 1 に記載の方法。

【請求項 6】

前記バイナリ・データ・ベクトルの前記アクセスが読み取り動作である、請求項 1 に記載の方法。

【請求項 7】

前記バイナリ・データ・ベクトルの前記アクセスが書き込み動作である、請求項 1 に記載の方法。

30

【請求項 8】

前記開始アドレスの前記選択された $\log_2(n)$ 個の部分が同じサイズに設定される、請求項 2 に記載の方法。

【請求項 9】

前記バイナリ・データ・ベクトルが、前記複数の前記 n 個のメモリ・バンク内の同じサイズの部分に格納される、請求項 1 に記載の方法。

【請求項 10】

前記開始アドレスの前記提供、前記オフセットの前記決定、ならびに前記バイナリ・データ・ベクトルの部分のシーケンスの前記決定、および部分のシーケンスの前記決定が、前記 n 個のメモリ・バンクへの前記バイナリ・データ・ベクトルの前記書き込み動作の前に実行される、請求項 7 に記載の方法。

40

【請求項 11】

前記開始アドレスの前記提供、前記オフセットの前記決定、および前記バイナリ・データ・ベクトルの部分のシーケンスの前記決定が、前記 n 個のメモリ・バンクに対する前記バイナリ・データ・ベクトルの前記読み取り動作の前に実行され、部分のシーケンスの前記決定が、前記 n 個のメモリ・バンクに対する前記バイナリ・データ・ベクトルの前記読み取り動作の後に実行される、請求項 6 に記載の方法。

【請求項 12】

バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んで

50

いるメモリ・ユニット内の前記バイナリ・データ・ベクトルにアクセスするためのメモリ・アクセス・ユニットであって、前記メモリ・アクセス・ユニットが、

前記バイナリ・データ・ベクトルの開始アドレスおよび前記バイナリ・データ・ベクトルの要素の2の累乗のストライド s を受信するように適応された受信ユニットと、

前記メモリ・バンクの各々について1つの n 個のオフセットを決定するように適応された第1の決定モジュールであって、前記オフセットの各々が、

複数のビットレベルのXOR関数を前記開始アドレスに適用してZベクトルを生成することと、

マッピング・テーブルにアクセスするために前記Zベクトルを使用することと、

前記バイナリ・データ・ベクトルの2の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることとによって決定される、前記第1の決定モジュールと、

前記Zベクトルの2進等価値に応じて前記 n 個のメモリ・バンク内の前記バイナリ・データ・ベクトルの部分のシーケンスを決定するように適応された第2の決定モジュールと、

前記メモリ・ユニットの前記 n 個のメモリ・バンク内の前記バイナリ・データ・ベクトルに並列にアクセスするように適応された基本アクセス・ユニットとを備える、メモリ・アクセス・ユニット。

【請求項13】

前記第1の決定モジュールによって、複数のビットレベルのXOR関数を前記開始アドレスに前記適用することが、

前記開始アドレスから $\log_2(n)$ 個の部分および $\log_2(n)$ ビット・サイズのYの部分を選択することと、

前記 $\log_2(n)$ 個の部分の各々についてパリティ・ビットを決定し、 $\log_2(n)$ 次元のパリティ・ベクトルXを生成することと、

XとYの間でビットレベルのXOR演算を実行し、Zベクトルを生成することを含む、請求項12に記載のメモリ・アクセス・ユニット。

【請求項14】

前記第1の決定モジュールによって、マッピング・テーブルにアクセスするために前記Zベクトルを使用することが、

前記マッピング・テーブル内のインデックスとして前記Zベクトルを使用することによって、 $n * n$ 個のエントリを含んでいる前記マッピング・テーブル内の行を選択することと、

前記マッピング・テーブルのアクセス結果を前記開始アドレスの部分と結合して、前記 n 個のオフセットを取得することとも含む、請求項13に記載のメモリ・アクセス・ユニット。

【請求項15】

前記第1の決定モジュールの前記結合が、

前記開始アドレスのビット・オフセット $\log_2(n)$ で開始する $\log_2(n)$ ビットを、前記決定された $\log_2(n)$ 個のマッピング・テーブル・エントリに置き換え、 n 個のオフセットを生成することも含み、各オフセットが、前記 n 個のバンクの各1つにおいてオフセットとして使用される、請求項14に記載のメモリ・アクセス・ユニット。

【請求項16】

前記第2の決定モジュールが、

前記 n 個のメモリ・バンク内の前記バイナリ・データ・ベクトルの前記部分のシーケンスを前記決定するときに、

$\log_2(n)$ 個の連続するマルチプレクサ段を制御するようにも適応され、前記マルチプレクサ段の各々が、前記バイナリ・データ・ベクトルの全幅をカバーするように、複数の2入力ビット・マルチプレクサを含む、請求項12に記載のメモリ・アクセス・ユニット。

【請求項17】

前記バイナリ・データ・ベクトルの前記アクセスが読み取り動作である、請求項12に

10

20

30

40

50

記載のメモリ・アクセス・ユニット。

【請求項 18】

前記バイナリ・データ・ベクトルの前記アクセスが書き込み動作である、請求項 12 に記載のメモリ・アクセス・ユニット。

【請求項 19】

前記開始アドレスの前記選択された $\log_2(n)$ 個の部分が同じサイズに設定される、請求項 13 に記載のメモリ・アクセス・ユニット。

【請求項 20】

前記バイナリ・データ・ベクトルが、前記複数の前記 n 個のメモリ・バンク内の同じサイズの部分に格納される、請求項 12 に記載のメモリ・アクセス・ユニット。

10

【請求項 21】

前記開始アドレスの前記提供、前記オフセットの前記決定、ならびに前記バイナリ・データ・ベクトルの部分のシーケンスの前記決定、および部分のシーケンスの前記決定が、前記 n 個のメモリ・バンクへの前記バイナリ・データ・ベクトルの前記書き込み動作の前に実行される、請求項 18 に記載のメモリ・アクセス・ユニット。

【請求項 22】

前記開始アドレスの前記提供、前記オフセットの前記決定、および前記バイナリ・データ・ベクトルの部分のシーケンスの前記決定が、前記 n 個のメモリ・バンクに対する前記バイナリ・データ・ベクトルの前記読み取り動作の前に実行され、部分のシーケンスの前記決定が、前記 n 個のメモリ・バンクに対する前記バイナリ・データ・ベクトルの前記読み取り動作の後に実行される、請求項 17 に記載のメモリ・アクセス・ユニット。

20

【請求項 23】

コンピュータ・プログラムであって、請求項 1 ないし 11 のいずれか 1 項に記載の方法の各ステップをコンピュータに実行させるための、コンピュータ・プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本明細書に記載された実施形態例は、一般に、効率的なメモリ・アクセスに関連しており、より詳細には、バイナリ・データ・ベクトルが一部に格納される複数のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための方法に関連している。実施形態例は、バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための関連するメモリ・アクセス・ユニットと、コンピュータ・プログラム製品とに、さらに関連する。

30

【背景技術】

【0002】

多くのワークロードおよび最新のコンピューティングは、同一の動作を適用する基本的な処理ステップに依存することに基づき、大量のデータに対する一連の動作である。そのようなワークロードで高性能を実現するための主な課題は、通常、最も効率的な方法でメモリ・システムから実行ユニットにデータを取得すること、および処理結果を実行ユニットからメモリ・システムに返すことに関連している。

40

【0003】

そのようなワークロードの多くによって必要とされる処理は、多くの場合、浮動小数点値または整数値のいずれかに対する乗算および加算の組み合わせである。例えば、深層学習の中心的処理（トレーニングおよび干渉の両方）のワークロードは、膨大な量の行列乗算動作を入力データまたはトレーニング・データあるいはその両方に適用することに基づく。ステンシル・コード（stencil codes）は、気象予測、石油およびガスの検索などのアプリケーションに使用される多くのモデルの中核部を形成する。これらのステンシル・コードは、隣接する要素の値の何らかの加重和を通常は含む固定パターン（ステンシルと呼ばれる）に従って、誤差率要素（error rate elements）を反復的に更新する。FFT

50

(Fast Fourier Transformation : 高速フーリエ変換) は、科学および工学において最も広く使用されているアルゴリズムの1つであり、非常に多くのいわゆるバタフライ演算を計算することに基づき、各バタフライ演算は、乗算および加算に基づく。ステンシル計算およびFFTの両方は、特定の深層学習アルゴリズム(例えば、畳み込みニューラルネットワーク(CNN: convolutional neural networks))にも使用される。

【0004】

高性能を獲得するために、実行ユニット(通常は、乗算器および加算器の組み合わせを含む)は、多くの場合、広いデータ・ベクトルのストリームとして提供される入力データに対して並列に、同じ動作または動作の組み合わせを何度も適用するベクトル(SIMD、単一の命令、複数のデータ)ユニットとして編成される。

10

【0005】

ベクトル・ユニットを構成する処理要素のためのすべての必要なオペランド値を含んでいるメモリ・システムからデータ・ベクトルを取り出すことは、多くの場合、それらのオペランド値を読み取るために必要なアドレス・シーケンス内のいわゆる2の累乗のストライドのため、簡単ではない。この種類のアクセス・パターンは、例えば、1のインクリメント(例えば、0、1、2、3など)、8のインクリメント(例えば、0、8、16、24など)、または256のインクリメント(例えば、0、256、512、768、1024など)などを伴うアドレスによって、線形仮想アドレス空間または物理アドレス空間上にマッピングされたデータ構造にアクセスするワークロードに発生する。

【0006】

バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための方法に関連する複数の開示が存在する。

20

【0007】

文書US 9705532 B2は、低密度パリティチェック・コードの生成の一部としての情報ビットの並列累算のための方法およびシステムを開示する。それによって、ベクトル演算を介して、連続する情報ビットが累算されることができ、累算に使用されるパリティ・アドレスは、仮想プライベート・パリティ・アドレス・マッピング(virtual and private parity address mapping)を介して連続的にされ得る。

【0008】

文書US 7793038 B2は、メモリ・バンク識別子およびメモリ・バンクの内部アドレスを含んでいる組み合わせアドレスへの、メモリ・システム内のメモリ・アクセスのためのアドレスのマッピングを開示する。このアドレスは、第1の部分および第2の部分に分割される。さらに、メモリ・バンク識別子は、検索行列内で検索動作を実行することによって決定される。

30

【0009】

しかし、多くのアプリケーションでは、データ・ベクトル要素にアクセスするために使用されるアドレス・シーケンス内の複数の異なる2の累乗のストライドを含むさまざまな方法でデータ・ベクトルを構築する必要があるという、既知の解決策の不利益が引き続き問題になる。例えば、あるときには、データ・ベクトルは、ストライド1を伴うアドレス・シーケンスに関連するデータ要素で構成されることがあり、別のときには、データ・ベクトルは、ストライド8を伴うアドレス・シーケンスに関連するデータ要素で構成されることがある、などである。同じデータ要素が、さまざまな方法で他のデータ要素と組み合わせられてデータ・ベクトルを形成するため、1つのデータ・ベクトルを構成するすべてのデータ要素を、基礎になるハードウェアの1原理サイクル(principle cycle)でそれらのバンクから並列に読み取ることができるよう、データ・ベクトルを形成するために使用されるすべてのデータ要素が個別のバンクに格納されることを保証することが、よりいっそう困難になる。本明細書において提案される概念は、この問題に対処する。

40

【発明の概要】

【0010】

50

1つの態様によれば、バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための方法が提供されてよい。この方法は、バイナリ・データ・ベクトルの開始アドレスおよびバイナリ・データ・ベクトルの要素の2の累乗のストライドsを受信することと、メモリ・バンクの各々について1つのn個のオフセットを決定することを含み、オフセットの各々は、複数のビットレベルのXOR関数を開始アドレスに適用してZベクトルを生成し、マッピング・テーブルにアクセスするためにZベクトルを使用し、バイナリ・データ・ベクトルの2の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることによって決定される。

【0011】

さらに、この方法は、Zベクトルの2進等価値(binary equivalent value)に応じてn個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定することと、メモリ・ユニットのn個のメモリ・バンク内のバイナリ・データ・ベクトルに並列にアクセスすることを含む。

【0012】

別の態様によれば、バイナリ・データ・ベクトルが一部に格納されてよい複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするためのメモリ・アクセス・ユニットが提供されてよい。メモリ・アクセス・ユニットは、バイナリ・データ・ベクトルの開始アドレスおよびバイナリ・データ・ベクトルの要素の2の累乗のストライドsを受信するように適応された受信ユニットと、メモリ・バンクの各々について1つのn個のオフセットを決定するように適応された第1の決定モジュールとを備え、オフセットの各々は、複数のビットレベルのXOR関数を開始アドレスに適用してZベクトルを生成し、マッピング・テーブルにアクセスするためにZベクトルを使用し、バイナリ・データ・ベクトルの2の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることによって決定される。

【0013】

さらに、メモリ・アクセス・ユニットは、Zベクトルの2進等価値に応じてn個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定するように適応された第2の決定モジュールと、メモリ・ユニットのn個のメモリ・バンク内のバイナリ・データ・ベクトルに並列にアクセスするように適応された基本アクセス・ユニットとを備える。

【0014】

バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための提案された方法は、次のような複数の利点および技術的効果提供することができる。

【0015】

提案された方法は、特に、乗算などのような行列演算の場合に、「1回で」、すなわち、基本的にメモリ・システムの1動作サイクル内で、大きいデータ・ベクトルまたは長いデータ・ベクトルに対処することも、主として可能にすることができる。これは、重複するアドレス指定が必要とされなくてよいため、可能になることができる。高次元ベクトルの異なる次元のデータ値は、メモリ・システムまたはメモリ・ユニットの異なるメモリ・バンクにわたって常に分散されてよい。したがって、異なるメモリ・バンク、およびメモリ・バンク内の高次元ベクトルの1つの次元の1つの特定のデータ値をアドレス指定することによって、ベクトルのすべての値への最短のアクセス時間が可能になる。

【0016】

したがって、同じ超次元ベクトルの異なる次元の異なるデータ値をアドレス指定するために、同じメモリ・バンクへの2回以上のアクセスが必要とされることがある状況が回避されることができる。これによって、メモリ・ユニット、およびしたがって、関連するコンピューティング・システムの動作サイクルを明確に節約することができる。したがって、本明細書で提案される概念によって、従来技術と比較した場合に、メモリ・システム/

10

20

30

40

50

ユニット内でデータ・ベクトルにアクセスするための性能上の大幅な優位性を実現することができる。

【0017】

提案された概念は、メモリ・バンクに書き込むときにベクトルのデータ要素をシャッフルする方法、およびベクトルのデータを読み取るときにデータ要素を逆シャッフルする方法も、可能にする。したがって、書き込み・読み取りサイクルの後に、ベクトルのデータ要素は、メモリ・ユニット/システムへの書き込み動作の前のように、再び適切に格納される。

【0018】

さらに、コンピュータ・システム設計における現在の傾向に起因して、メモリ内のデータへのアクセスおよび対応する転送が、多くのアプリケーションのシステムレベルの性能に影響を与える支配的要因になる。最先端のベクトル・ユニットによって極めて効率的に処理され得るデータ・ベクトルにメモリ内でアクセスできる速度を改善することによって、本明細書に記載された実施形態は、ある範囲のワークロードに対するコンピュータ・システムの性能を大幅に改善するのに役立つことができる。

10

【0019】

以下では、本発明の概念の追加の実施形態が説明される。

【0020】

この方法の1つの好ましい実施形態によれば、複数のビットレベルのXOR関数を開始アドレスに適用することは、開始アドレス(具体的には、CPUの線形アドレス空間内のベクトルの第1の要素のアドレス)から、 $\log_2(n)$ 個の部分および $\log_2(n)$ ビット・サイズのYの部分を選択することと、 $\log_2(n)$ 個の部分の各々についてパリティ・ビットを決定し、 $\log_2(n)$ 次元のパリティ・ベクトルX(すなわち、個別のビットのグループ)を生成することと、XとYの間でビットレベルのXOR演算を実行してZベクトルを生成することとを含んでよい。それによって、アドレス・マスクによって(具体的には、アドレス・マスクによる $\log_2(n)+1$]によって) $\log_2(n)$ 個の部分およびYの部分を選択可能であってよいということが述べられてよい。このシャッフル動作の詳細は、下で図を説明するとき、組み合わせにおいてさらに明確になるであろう。

20

【0021】

方法の1つの有用な実施形態によれば、マッピング・テーブルにアクセスするためにZベクトルを使用することは、マッピング・テーブル内のインデックスとしてZベクトル(具体的には、その2進等価値)を使用することによって、 $\log_2(n)$ 個のエントリを含んでいるマッピング・テーブル内の行を選択することと、マッピング・テーブルのアクセス結果を開始アドレスの部分と結合して(「マージする」と表すこともできる)n個のオフセットを取得することとを含んでもよい。下で、図との関連において詳細が説明される。

30

【0022】

方法の1つの有利な実施形態によれば、この結合することは、開始アドレスのビット・オフセット $\log_2(s)$ にある $\log_2(n)$ ビットを、決定された $\log_2(n)$ 個のマッピング・テーブル・エントリに置き換え、 $\log_2(n)$ 個のオフセットを生成することを含んでよい。それによって、各オフセットが、 $\log_2(n)$ 個のバンクの各1つにおいてオフセットとして使用されてよい。これによって、極めて限定されたオーバーヘッドのみを伴う簡単な決定を可能にすることができる。

40

【0023】

方法の1つのさらに好ましい実施形態によれば、n個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定すること(シャッフルとして表されてもよいプロセス)は、 $\log_2(n)$ 個の連続するマルチプレクサ段を制御することを含んでよく、マルチプレクサ段の各々は、データ・ベクトルの全幅をカバーするように、複数の2入力ビット・マルチプレクサ(2-input-bit-multiplexers)を含んでいる。これによって

50

、重複するアドレスが必要とされなくてよいように、高次元ベクトルの個別のデータ値の位置決めを可能にすることができる。

【0024】

より詳細には、対応するZベクトルのビットが1に等しい場合に、第1のマルチプレクサ段が、データ・ベクトルの最上位部分を、最下位部分と交換してよい。第2のマルチプレクサ段が、第2のZベクトルのビットの制御下で、データ・ベクトルの下位部分および上位部分に対して、同じ機能を並列に適用する、などであってよい。これらの $\log_2(n)$ 個のマルチプレクサ段は、メモリ・バンクの正しい位置に書き込まれ、メモリ・バンクの正しい位置からそれぞれ読み取られる、さまざまなデータ要素を交換するために使用されてよい。

10

【0025】

方法の1つの許容される実施形態によれば、アクセス動作が読み取り動作であってよい。代替として、または組み合わせて、アクセス動作が書き込み動作であってよい。したがって、本発明の概念の一部として、削除操作も使用可能である。したがって、個別のメモリ・バンクのメモリ・ユニット、メモリ・デバイス内のすべての動作が、本明細書で提案される概念を使用して実行されてよい。より大きい(例えば、超次元)ベクトルの識別されたデータ要素のアドレス生成または再シャッフルのシーケンスには、差異が存在してよい。

【0026】

方法の1つの実用的な実施形態によれば、開始アドレスの選択された $\log_2(n)$ 個の部分は、同じサイズに設定されてよい。これは、2の累乗のストライドの概念の下で同様に実行されてよい。

20

【0027】

方法の別の実用的な実施形態によれば、バイナリ・データ・ベクトルが、複数のn個のメモリ・バンク内の同じサイズの部分に格納されてよい。したがって、データ要素の等しい分散が結果として実現されてよい。

【0028】

メモリ・バンクへの書き込み動作が対象にされる方法の1つの実施形態によれば、開始アドレスを提供することと、オフセットを決定することと、バイナリ・データ・ベクトルの部分のシーケンスを決定することと(すなわち、主概念の部分)、部分のシーケンスを決定することと(すなわち、シャッフル)が、n個のメモリ・バンクへのバイナリ・データ・ベクトルの書き込み動作の前に実行されてよい。これによって、メモリ・ユニットの1サイクルのみを基本的に必要とする高速書き込み動作のため、およびデータに対する「ミラーリングされた」読み取り動作を可能にするために、ベクトルの次元値を格納することを保証することができる。

30

【0029】

方法の1つの結果とし生じる実施形態によれば、開始アドレスを提供することと、オフセットを決定することと、バイナリ・データ・ベクトルの部分のシーケンスを決定することと(主概念を参照)は、n個のメモリ・バンクに対するバイナリ・データ・ベクトルの読み取り動作の前に実行されてよく、部分のシーケンスを決定すること(すなわち、シャッフル)は、n個のメモリ・バンクに対するバイナリ・データ・ベクトルの読み取り動作の後に実行されてよい。このステップは、大きい次元のベクトルまたは超次元ベクトルのデータ値を元のシーケンスに戻すために必要になることがある。

40

【0030】

さらに、実施形態は、コンピュータまたは任意の命令実行システムによって、またはこれらに接続して使用するためのプログラム・コードを提供するコンピュータ使用可能媒体またはコンピュータ可読媒体からアクセスできる関連するコンピュータ・プログラム製品の形態を取ってよい。この説明の目的で、コンピュータ使用可能媒体またはコンピュータ可読媒体は、命令実行システム、命令実行装置、または命令実行デバイスによって、またはこれらに接続して使用するためのプログラムを格納するか、伝達するか、伝搬するか、

50

または運ぶための手段を含むことができる任意の装置であってよい。

【0031】

本発明の実施形態が、さまざまな対象を参照して説明されるということに注意する必要がある。具体的には、一部の実施形態は、方法タイプの請求項を参照して説明され、他の実施形態は、装置タイプの請求項を参照して説明される。ただし、当業者は、前述の説明および以下の説明から、特に注記のない限り、対象の1つの種類に属している特徴の任意の組み合わせに加えて、異なる対象に関連する特徴間、具体的には、方法タイプの請求項の特徴と、装置タイプの請求項の特徴との間の任意の組み合わせも、本文書内で開示されると見なされるということ推測するであろう。

【0032】

上で定義された態様およびその他の態様は、以下で説明される実施形態の例から明らかになり、実施形態の例を参照して説明されるが、実施形態はこれらに限定されない。

【0033】

以下の図面を単なる例として参照し、好ましい実施形態について説明する。

【図面の簡単な説明】

【0034】

【図1】バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための本発明の方法の基本実施形態のブロック図を示す図である。

【図2】標準的な行列変更のために必要な困難なメモリ・アドレス指定方式を示す図である。

【図3】提案された概念の別の応用分野を示す図である。

【図4】例示的なマッピングの実施形態を示す図である。

【図5】9ビットの要素を含んでいるデータ・ベクトルの開始アドレスの例を示す図である。

【図6】例示的なマッピング・テーブルおよびメモリ・ユニット内のメモリ・バンクのオフセットを決定することにおける次のステップの実施形態を示す図である。

【図7】メモリ・バンク内のエントリのシャッフルの決定方法を示す図である。

【図8】メモリ・ユニットを表す例示的な8つのバンクを示す図である。

【図9】ベクトル開始アドレス256、ストライド1としての例を示す図である。

【図10】ベクトル開始アドレス3、ストライド8としての例を示す図である。

【図11】提案された基本概念に従って、より実装に近い読み取りプロセスのフローチャートを示す図である。

【図12】提案された基本概念に従って、より実装に近い読み取りプロセスのフローチャートを示す図である。

【図13】複数のバンクを含んでいるメモリ・ユニット内のバイナリ・ベクトルのためのメモリ・アクセス・ユニットのブロック図を示す図である。

【図14】メモリ・アクセス・ユニットが使用され得るコンピュータ・システムのブロック図を示す図である。

【発明を実施するための形態】

【0035】

本説明の文脈において、以下の規則、用語、または表現、あるいはその組み合わせが使用されてよい。

【0036】

「バイナリ・データ・ベクトル」という用語は、コンピュータのメモリ・システム内のビットのシーケンスのワードとしてアドレス指定される値のフィールド（例えば、ビット値、整数値、浮動小数点値、文字フィールドなど）を示してよい。そのようなメモリ・システムは、通常、RAM（random access memory：ランダム・アクセス・メモリ）であってよい。特定のアドレスでのワード長に応じて、事前に定義された数のビットがメモリ・システムまたはメモリ・ユニット（テキスト全体を通じて、両方の形態が同義語とし

10

20

30

40

50

て使用されてよい)のワードに格納されてよい。

【0037】

「メモリ・ユニット」という用語は、複数のバンクに分割されることがあるコンピュータ・システムのより大きいストレージ・モジュールを示してよい。各バンクは、特定のサイズの個別にアドレス指定可能なワードに編成された、事前に定義された数のストレージ・セルを含んでよい。メモリ・ユニットのバンクのメモリ・バンクは、個別にアドレス指定可能であってよく、各メモリ・バンク内で、個別のワードがアドレス指定可能であってよい。

【0038】

「メモリ・バンク」という用語は、ストレージ・セルのグループを示してよく、この用語の複数形は、バンク・アドレスを介してアクセスできるストレージ・セルのより大きいグループを定義する。各バンク内で、ワードとして編成された複数の個別のメモリ・セルが、個別にアドレス指定されてよい。通常、メモリ・ユニットは、複数のメモリ・バンクを含む。

10

【0039】

「バイナリ・データ・ベクトルの開始アドレス」という用語は、通常、長いバイナリ・データ・ベクトルの第1の(または第0の)要素のアドレスを示してよい。

【0040】

「2の累乗のストライド」という用語は、 2^n 個の要素によって間を分離されている要素のシーケンスを示してよい。 $n = 0$ である場合、 $2^0 = 1$ であるため、要素0、1、2、3、...がアドレス指定されてよい。 $n = 1$ である場合、 $2^1 = 2$ であるため、要素0、2、4、6、...がアドレス指定されてよい、などとなる。

20

【0041】

「Zベクトル」という用語は、複数のビットレベルのXOR関数をアドレスに適用することによって取得された $\log_2(n)$ ビットのサイズを有するベクトルを示してよい。

【0042】

「マッピング・テーブル」という用語は、メモリ・ユニットに含まれているバンクの数と同じ数の列を含んでいるテーブルを示してよい。テーブル内の要素は、0からバンクの数 - 1までの整数値であってよい。

【0043】

以下では、各図について詳細に説明する。図内のすべての命令は概略図である。まず、バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための本発明の方法の実施形態のブロック図が示される。その後、さらに別の実施形態に加えて、バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするためのメモリ・アクセス・ユニットの実施形態が説明される。

30

【0044】

図1は、バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための方法の実施形態のブロック図を示している。以下の例の精神的参照モデルとして、バイナリ・データ・ベクトルが256ビットの長さであること、およびメモリ・ユニットが8つのバンクを含むこと、各バンクの32ビットの幅、ならびに64個の位置のバンク・サイズを仮定することができる。

40

【0045】

この方法は、バイナリ・データ・ベクトルの開始アドレス(具体的には、CPUのアドレス空間内の配列の開始アドレス)およびバイナリ・データ・ベクトルの要素の2の累乗のストライド s を受信すること(102)を含む。

【0046】

この方法は、メモリ・バンクの各々について1つの n 個のオフセットを決定すること(

50

104) も含む。それによって、複数のビットレベルの XOR 関数を開始アドレスに適用して Z ベクトルを生成することと、マッピング・テーブルにアクセスするために Z ベクトルを使用することと (108)、バイナリ・データ・ベクトルの 2 の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることと (110) によって、オフセットの各々が決定される (106)。

【0047】

さらに、この方法は、Z ベクトルの 2 進等価値に応じて n 個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定することと (112)、メモリ・ユニットの n 個のメモリ・バンク内のバイナリ・データ・ベクトルに並列に (具体的には、1 動作サイクル内で、または同時に) アクセスすることと (114) を含む。

10

【0048】

図 2 は、200 に、標準的な行列演算のために必要な困難なメモリ・アドレス指定方式を示している。この図は、単純な行列乗算の例を示しており、行列 A の行が、要素ごとに行列 B の列に適用される。両方の行列の要素が、線形アドレス空間内の連続するアドレス位置に書き込まれる場合、行列 A 内の各行は、アドレス・ストライド 1 を使用してアクセスされ、行列 B 内の各列は、行のサイズに等しいアドレス・ストライドを使用してアクセスされる。多くのメモリ・システムでは、後者のストライドが問題になる可能性があるため、通常、「行列の対角線上」の行列値を反転することを伴う、いわゆる行列転置演算が適用される。その結果、次に、ストライド 1 のアクセス・パターンを使用して列にアクセスし、結果として生じる行列 C を生成することができる。しかし、これは、行列全体を読み取って再び書き込む必要があるという犠牲を伴い、大量のメモリ帯域幅および時間を消費する可能性もある。

20

【0049】

提案された概念の別の応用分野を示す別の例が、図 3 の 300 に示されている。図 3 は、2D 構造または 3D 構造に適用できるさまざまなステンシルを示している。この図から、単一のステンシルによってカバーされたすべての要素にアクセスすることが、(2D 構造および 3D 構造の場合に) それぞれ 2 つ、3 つの異なる 2 の累乗のストライドを含んでいるということがわかる。

【0050】

長い間、2 の累乗のアクセス・パターンが調査されたにもかかわらず、非効率的アクセスの問題はまだ解決されていない。1 つの問題は、多くの場合、2 の累乗の数のメモリ・バンクを使用して編成された多くのメモリ・システムにおいて使用される編成および卓越したアドレス・マッピングに起因して、このような種類のアクセスが、メモリ性能の低下をもたらすことが多いということである。

30

【0051】

特殊な HPC (high performance computing: 高性能コンピューティング) システムでは、素数に対してモジュロ演算を実行すること、または同様の方法に基づく固定された / 配線接続されたマッピングが適用されることがある。これらのマッピングは、異なっているか、またはプログラムすることが不可能であるため、特定のワークロードに固有であることがあり、したがって、いまだに悪い状況 / 異常な状況にある可能性がある。後者のメカニズムも、アクセス待ち時間を増加させる。

40

【0052】

次の図に進む前に、ベクトルの処理およびアクセスに関する次の考察について検討すべきである。

【0053】

n 個のバンク (すなわち、 $n * w$ に等しいデータ・ベクトルのサイズであり、w は単一のバンクの幅またはアクセス粒度である) 内のデータ・ベクトルに並列にアクセスするために、バンクごとに正しいオフセットが生成され、データ・ベクトルごとに合計 n 個のオフセットを得る必要がある。例えば、図 4 のマッピング例 400 に基づいてアドレス・シーケンス 0、1、2、3、4、5、6、および 7 (ストライド 1) に関連する 8 要素のデ

50

ータ・ベクトルにアクセスする必要がある場合、8つの各メモリ・バンク内で、次のオフセットでメモリ位置にアクセスする必要がある： $\{0, 1, 2, 3, 4, 5, 6, 7\}$ 。

【0054】

示されているテーブルは、複数のメモリ・バンクから成るメモリ・システムに格納された1つのデータ構造への2の累乗のストライドのアクセス・パターン間で、完全にインターリーブすることができる、柔軟でプログラム可能なアドレス・マッピング方式を定義することによって、前述した問題に対処する。

【0055】

このマッピング方式は、マッピングされるアドレスから抽出された2つの値によってインデックス付けされた小型の2Dルックアップ・テーブル（図示されていない）に基づく。2つのインデックス値のうちの一つは、選択されたアドレス・ビットのビット単位のパリティ値を計算することによって決定される。しかし、基礎になる概念についてここで説明する必要はない。この図は、完全なインターリーブの複数の2の累乗アクセス・パターンの場合、メモリ位置のシーケンスのシャッフルも必要になることがあるということのみを示しているはずである。

【0056】

図4の編模様でマーク付けされた「マッピングされたアドレス」は、2の累乗のストライドを伴うアライメントされたアドレス・シーケンスに関連する（例示的な8つのデータ要素で）メモリ・バンクの数に等しいサイズを有するデータ・ベクトルが、すべてのメモリ・バンク上で完全にインターリーブされており、したがって各バンク上で正確に1回のアクセスを必要として、すべてのバンクに並列にアクセスすることによって、単一のデータ・ベクトルとしてアクセスできる（すなわち、読み取るか、または書き込むことができる）ことを示している。例えば、図4では、ストライド1を伴うアライメントされたアドレス（例えば、アドレス0、1、2、3、4、5、6、7）、ストライド2を伴うアライメントされたアドレス（例えば、アドレス0、2、4、6、8、10、12、14）、ストライド4を伴うアライメントされたアドレス（例えば、アドレス0、4、8、12、16、20、24、28）、ストライド8を伴うアライメントされたアドレス（例えば、アドレス0、8、16、24、32、40、48、56）、ストライド16を伴うアライメントされたアドレス（例えば、アドレス0、16、32、48、64、80、96、112）、ストライド32を伴うアライメントされたアドレス（例えば、アドレス0、32、64、96、128、160、192、224）などの各ブロックに、すべて並列にアクセスすることができる。アクセスされるデータ要素がバンク間で均等に分散されるため、バンクの再アクセスは不要である。

【0057】

読み取り動作の場合、バンク0がデータ・ベクトルの第1の要素（マッピングされたアドレス0）を提供し、バンク1がデータ・ベクトルの第2の要素（マッピングされたアドレス1）を提供する、などとなる。同様に、書き込み動作の場合、データ・ベクトルの要素が、連続的順序で各バンクに書き込まれる。代わりに、アドレス・シーケンス0、2、4、6、8、10、12、および14（ストライド2）に関連する8要素のデータ・ベクトルがアクセスされる場合、8つの各バンク内で、次のオフセットに並列にアクセスする必要がある： $\{0, 8, 2, 10, 4, 12, 6, 14\}$ 。この例では、読み取り動作の場合、データ・ベクトルの第1の要素（マッピングされたアドレス0）がバンク0によって提供され、データ・ベクトルの第2の要素（マッピングされたアドレス2）がバンク2によって提供され、データ・ベクトルの第5の要素（マッピングされたアドレス8）がバンク1によって提供される、などとなる。マッピングされたアドレスの順序（0、2、...、14）で要素を含んでいるべきである望ましいデータ・ベクトルを復元するために、バンクから読み取られたデータがシャッフルされなければならないということは明らかであろう。さらに、データ・ベクトルが書き込まれる場合、メモリ・バンク上で対応するアドレスがマッピングされる方法に従って各要素を正しく書き込むために、要

10

20

30

40

50

素もシャッフルされなければならない。

【 0 0 5 8 】

データ・ベクトルの要素に関連するアドレス・シーケンスに応じて、ある範囲の異なる順序が存在する可能性があり、その場合、個別のバンクのオフセットが生成される必要があり、読み取りデータおよび書き込みデータがシャッフルされる必要がある。

【 0 0 5 9 】

総当たりの方法は、他の要素とは無関係にデータ・ベクトル内の要素ごとにアドレスを生成し、バンク識別子およびバンク・オフセットに独立してマッピングする方法であり、合計 n 個のアドレスの生成機能およびアドレス・マッピング機能をもたらす。次に、マッピング機能の出力（具体的には、バンク識別子）は、データ・ベクトルを書き込むためおよび読み取るための 2 つの n 対 n シャッフル・ユニットを制御し、任意の可能な順序での各データ・ベクトルの要素のシャッフルを可能にするために使用され得る。

10

【 0 0 6 0 】

しかし、性能および面積の両方のコストに関して、 n の値と共に、この方法の複雑さが急激に増大し、総当たりの方法を、 n の小さい値（例えば、4 または 8）に適用することに基本的に制限する。

【 0 0 6 1 】

しかし、より大きいデータ・ベクトルは、計算における必要な効率および性能を達成するのに役立つため、非常に重要である。このことは、より新しいメモリ技術（例えば、HBM (high bandwidth memory: 高帯域幅メモリ)）によって提供される広いデータ・パスのために、よりいっそう強調される。

20

【 0 0 6 2 】

そのため、次の技術革新が対象とされる。

1. 並列にすべてのバンクをカバーするデータ・ベクトル全体のベクトルに基づくオフセット生成。
2. メモリ・バンクに書き込むためのデータ・ベクトル要素の効率的なシャッフル。
3. メモリ・バンクから読み取った後のデータ・ベクトル要素の効率的なシャッフル。
4. パリティ・ビットの各々の計算に使用されるアドレス・ビットの選択。

【 0 0 6 3 】

ベクトルに基づくオフセット生成の場合、任意の 2 の累乗のストライドを伴うアライメントされたアドレス・シーケンスに関連する n 個の要素を含む任意のデータ・ベクトルに関して、 n 個のアドレスが、アドレスにおける $\log_2 n$ 個の隣接するビット位置に制限されることと、それらの $\log_2 n$ 個のビット位置の n 個の可能な値をすべてカバーすることとの間の違いについて、以下の観察が行われ得る。

30

【 0 0 6 4 】

ビット・インターリーブ方式でアドレスからパリティ・ビットが抽出される場合、2 の累乗のストライドを伴う任意のアライメントされたアドレス・シーケンスに含まれる n 個のアドレスは、ルックアップ・テーブル内の単一の行で指定されたバンク識別子に基づいてマッピングされる。

【 0 0 6 5 】

図 5 は、9 ビットの要素 a_0 、 a_1 、 \dots 、 a_8 を含んでいるデータ・ベクトルの開始アドレス「a」502 の例 500 を示している。MSB (most significant bits: 最上位ビット) は、後で必要とされるベクトル Y として使用され、 $\log_2 (n)$ のサイズ ($n = 8$) である。ビット a_0 、 a_1 、 \dots 、 a_5 は、マスク・ベクトル m_0 、 m_1 、および m_2 ($\log_2 (n) = 3$) を使用して AND 結合され、図の右側の $\log_2 (n)$ 個のベクトル 504 を生成する。マスク・ベクトル m_0 、 m_1 、および m_2 が、開始アドレス 502 の Y の部分の位置に「0」を含んでいることに注意する。

40

【 0 0 6 6 】

ベクトル 504 の各々について、パリティ・ビットが計算され、 X_0 、 X_1 、 X_2 を生成する。これは、XOR 関数を適用することとして解釈されてもよい。 X_0 、 X_1 、 X_2 は、

50

Xベクトルの次元として扱われる。次に、ビットレベルのXOR関数がXおよびYに適用され、同じサイズまたは次元($\log_2(n)$)のベクトルZを生成する。

【0067】

これに関して、図6は、600で、メモリ・ユニット内のメモリ・バンクのオフセットを決定することにおける次のステップを示している。

【0068】

マッピング・テーブルが、X値およびY値の特定の組み合わせに関連して格納されている値が、それらのX値およびY値のXOR積に等しいという特性を有するということに注意する。これは、やはり、バンク識別子*i*が特定の行の列*j*に発生した場合、バンク識別子*j*がその行の列*i*にも発生する、という特性をもたらす(それらの値は、交換されて示され得る)。

10

【0069】

例えば、マッピング・テーブル602内のY=6に対応する行が、バンク識別子7を(X=1に対応する)列1に含んでいる。ここで、同じ行が、バンク識別子1を列7(X=7)に含んでいる。これら2つの位置は、「交換」位置の対と呼ばれる。上記の「交換」特性のため、*n*個のアドレスから抽出されたパリティ値が、*n*個のアドレスが異なっている、 $\log_2(n)$ 個の隣接するビット位置の値と同じ順序である場合、ここで示されるように、バンク・オフセットをルックアップ・テーブルの行から直接選択することができる。

【0070】

20

ベクトルZの2進等価値が、マッピング・テーブル602への行インデックスとして使用される。マッピング・テーブル602の行の各2進値が、開始アドレスの $\log_2(n)$ ビットを置き換える定義された位置で使用される。これが、図6の下側部分に示されている。アドレスにおいて(図5/502の部分a0、a1、...、a5を比較する)、ビット・オフセット3で開始する $\log_2(n)$ ビットが、各行エントリの2進等価値に置き換えられ、例えば、図6の左下部分からわかるように、「3」が011(2進数)を与える。この原理が、マッピング・テーブルの選択された行内の各要素に適用される。理解可能にするために、最初の2つおよび8番目のアドレス/行エントリの組み合わせのみが示されているということに注意する。この時点で、個別のバンクのオフセット値が決定されている。

30

【0071】

8ビットの例の場合、マッピング・テーブル(すなわち、ルックアップ・テーブル)が $8 * 8 * 3$ ビット=24バイトのみを必要とするということにも注意する。そのような非常に小さいテーブル内のエントリの検索は、提案された方法が表す極めて大きい優位性にとって、無視できる程度の要件しか表さない。

【0072】

次ステップでは(図7を比較する)、メモリ・バンク内のエントリのシャッフルを決定する必要がある。この決定は、Zベクトルの再呼び出しを必要とする。Zがバイナリ・ベクトルであるということに起因して、Zベクトルの要素は「0」または「1」のいずれかになる。したがって、ベクトルZのビットによって制御される交換制御は、ビット=1が開始アドレスでのエントリの対の交換を意味し、ビット=0が「交換なし」を意味する、というように定義される。

40

【0073】

これが、図7の下向き矢印(交換なし)または交差した矢印(交換)として示されており、これらの矢印は、図6に示されているように、ZベクトルのLSB(least significant bit: 最下位ビット)から次々に開始する $\log_2(n)$ ステップのシーケンスにおいて適用される。

【0074】

次にこれを、図8のメモリ・ユニット802を表す8つのバンクの例800において示すことができる。例えば、(スライド1を伴う)アドレス・シーケンス64、65、6

50

6、67、68、69、70、および71は、1に等しいY値を含み、したがって、マッピング・テーブル804内の行1(Y=1)に従ってマッピングされる。アドレス64がバンク1(806を比較する)にマッピングされ(808を比較する)、アドレス65がバンク0にマッピングされ、アドレス66がバンク3にマッピングされ、アドレス67がバンク2にマッピングされ、アドレス68がバンク5にマッピングされ、アドレス69がバンク4にマッピングされ、アドレス70がバンク7にマッピングされ、アドレス71がバンク6にマッピングされる。各オフセットは、6最下位アドレス・ビット0、1、2、3、4、5、6、および7に等しい。(前述したように、)これら8つのアドレスが異なっている $\log_2 8 = 3$ のビット位置は、アドレス・ビット位置0、1、および2にあり、これらのアドレス・ビット位置は、バンク・オフセットの3最下位ビットにも等しい。前の8つのオフセットの場合、これらのビット位置での値は、0、1、2、3、4、5、6、および7に等しい(これらは、この例では他のオフセット・ビットがたまたま0であるため、オフセット値に等しい)。ルックアップ・テーブルの内容が生成される方法のため、これらの異なる $\log_2(n)$ 個のビット位置の値は、パリティ・ビットの値に等しくなる。

10

【0075】

ルックアップ・テーブルの内容の上記の「交換」特性に基づいて、アドレスが、対応するルックアップ・テーブルの要素に基づいて特定のバンク識別子にマッピングされる場合、その要素に対応する交換位置が、異なる $\log_2(n)$ 個のビット位置の値に等しい、その同じアドレスのパリティ・ビットの値を含むということが導き出され得る。これは、アライメントされたアドレス・シーケンスにおいて第1のアドレスのY値を決定することによって、オフセット・ベクトルを生成することができ、その後、異なる $\log_2(n)$ 個のビット位置で、シーケンス内の第1のアドレスのY値に対応する各行の値を置換することによって、(前述したように、最下位部分を取得することによって)その第1のアドレスから取得されたオフセットのn個のコピーを取得することができるということを意味する。

20

【0076】

この特性は、任意の2の累乗のストライドを伴うn個のアドレスの任意のアライメントされたシーケンスに当てはまる。これは、前述したように、n個の異なるビット位置が、「アライメントされた」方法でXフィールドのビットをカバーする(すなわち、ビット位置X0から「開始」し、その後位置X1が続く、などとなる)、2の累乗のストライドに限定されない。これは、任意のn個の連続するビット位置の値が、常に、同じテーブルの行内の異なるエントリにインデックス付けする結果をもたらすからである。

30

【0077】

注目すべき示唆は、アドレス「a」から導出されたY値によってインデックス付けされた行からの値を、2の累乗のストライドの \log_2 の値に対応するn個のビット位置で、アドレスaから導出されたオフセット値に代入することによって、この特性に基づいて、n個のバンク・オフセットの値が直接取り出され得るということである。これが、次のステップを適用することによって、すべてのバンクのオフセットを効率的に生成するためのステップを示している、図11のフローチャートに示されている。(1)アドレス・シーケンスの第1のアドレスのみをマッピングし(すなわち、パリティを計算し、Y値を導出し、オフセットを導出し)、(2)Yに対応する1つのテーブルの検索を実行し、(3)並列なシフトおよびビット単位のOR関数を実行して、実際のバンク・オフセットを取得する。

40

【0078】

フローチャートを参照する前に、マッピングの2つの追加の例が与えられてよい。図9は、ベクトル開始アドレス256、ストライド1(アドレス256、257、258、. . .、263)としての例900を示している。256が100000000bに等しいため、マッピング・テーブル(図6の602を比較する)内の行4(100b)が4 5 6 7 0 1 2 3に等しくなるように、Yベクトル、Xベクトル、およびZベクトル

50

が、 $Y = 100b$ 、 $X = 000b$ 、 $Z = 100b$ であると決定され得る。したがって、(それぞれ、バンク0~7の)バンク・オフセットは、4、5、6、7、0、1、2、3である。 $Z = 100b$ を使用するデータ・ベクトルのシャッフルは、第1の4つのエントリが第2の4つのエントリと交換されるということの意味する。

【0079】

図10は、ベクトル開始アドレス3、ストライド8(したがって、アドレス3、11、19、...、59)としての例1000を示している。開始アドレス $= 3 = 000000011b$ であるため、 $Y = 000b$ 、 $X = 011b$ 、 $Z = 011b$ を得る。マッピング・テーブル602内の行3(011b)は、3 2 1 0 7 6 5 4に等しい。したがって、(それぞれ、バンク0~7の)バンク・オフセットは、27(011011b)、19(010011b)、11、3、59、51、43、35として決定され得る。

10

【0080】

データ・ベクトルのシャッフル部分は、次によって決定される。 $Z = 011b$ は、隣接するエントリの各対が、第1のステップで最初に交換され、それに続いて、2つのエントリの第1のブロックが第2のブロックと交換され、2つのエントリの第3のブロックが第4のブロックと交換されるということの意味する。その結果が、図8に、各バンクのエントリの横の縞模様マークで示されている。

【0081】

図11は、図1と比較した場合により実装に近い形態で、読み取りプロセスのアクセス・オフセット生成のフローチャート1100を示している。

20

【0082】

簡単に言うと、アクションのフローは次のように実行される(詳細が、さらに下でも説明される)。まず、アドレス「a」が取得される(1102)。次に、アドレス「a」のパリティ、Yフィールド、および初期オフセットが導出される(1104)。Yフィールド(またはパリティ)に基づいてマッピング・テーブル内の行にアクセスすることによって、すべてのバンクのベース・オフセットが取得される(1106)。次に、すべてのバンクのベース・オフセットが、対象の2の累乗のストライドに対応するビット位置にシフトされる(1108)。最後に、アドレス「a」に対して決定された初期オフセットが、マッピング・テーブルから取得されたシフトされたベース・オフセットとマージされ(1110)、すべてのメモリ・バンクのオフセットを取得する。

30

【0083】

この簡潔なフローチャート表現は、次の段落で再び説明される。

【0084】

データ・ベクトル要素をアドレス・シーケンス内のアドレスのメモリ・バンクに書き込むためのデータ・ベクトル要素の効率的なシャッフルは、 n 個のメモリ・バンクにわたって任意の可能な方法でマッピングすることができ、各データ要素を正しいバンクに書き込むために対応する書き込みデータをシャッフルする必要がある、合計 $n!$ 個の可能性が存在する。大きい値の n の場合、この可能性が急激に増大し、非常に大きい n 対 n シャッフル・ユニットが必要になる。

【0085】

40

しかし、前述したマッピング方式の場合、シャッフルを大幅に簡略化することができる。このことが、上記の例示的なマッピング方式を使用してここで説明される。

【0086】

以下では、パリティ値(すなわち、 X)が、オフセット値と完全に同じように順序付けられるということ仮定する。この仮定は、異なる最下位アドレス・ビットが最下位パリティ(X)ビットによってカバーされ、異なる第2のアドレス・ビットが次のパリティ・ビットによってカバーされる、などとなる場合に当てはまる。任意の2の累乗のストライドのアクセスで、この制約をどのように満たすことができるかについて、アドレス/パリティ・ビットの選択を対象にする以下の段落でさらに説明する。

【0087】

50

マッピング・テーブルが開始される方法（「交換」特性）のため、バンク・オフセット・ベクトルの生成と同様に、ルックアップ・テーブル内の単一の行によってシャッフルが決定される。次のようなマッピングの興味深い特性を、図6のマッピング・テーブルから直接理解することができる。マッピング・テーブル内の各行は、 k 個の連続する行要素をカバーするブロックを使用して、「アライメントされた方法」で行要素を交換することによって、（ $Y = 0$ に対応する）第1の行から導出されることができ、 k は2の累乗であり、2～行サイズ $n/2$ の範囲内の値を有する。

【0088】

図6の例では、2つの要素を含む4つのブロック（それぞれ、 $\{0, 1\}$ 、 $\{2, 3\}$ 、 $\{4, 5\}$ 、 $\{6, 7\}$ ）を識別し、各ブロック内すべての要素を交換することによって、 $Y = 1$ に対応する行が、 $Y = 0$ に対応する第1の行から導出され得る。 $Y = 2$ に対応する次の行は、最初に、4つの要素を含む2つのブロック（それぞれ、 $\{0, 1, 2, 3\}$ および $\{4, 5, 6, 7\}$ ）を識別し、次に、それらの各ブロック $\{0, 1\}$ および $\{2, 3\}$ 、 $\{4, 5\}$ および $\{6, 7\}$ の各々において2つの要素のサイズを有する2つのブロックを識別し、それらの2つの2要素ブロックの各々を4要素ブロック内で交換することによって導出され得る。4要素ブロックのレベルで、または4要素ブロック内の2要素ブロックのレベルで交換することによって、あるいはこれらの両方の組み合わせによって、残りのすべての行が第1の行から導出され得るということが、図6からわかる。

【0089】

図6および8の例からわかる別の重要な特性は、Zベクトルの最上位ビットを使用して、そのマッピング・テーブルの最大のブロック・サイズのレベルで交換動作を実行する必要があるかどうかを判定することができ、Zベクトルの次のビットを使用して、2番目に大きいブロック・サイズのレベルで交換動作を実行する必要があるかどうかを判定することができる、などとなるということである。

【0090】

「交換」特性のため、ここで、次のような方法で同じブロック交換の概念を用いて、前述した特性を、書き込みデータをシャッフルすることに直接使用できる。（1）（前述したように） n 個のアドレスのシーケンスの開始アドレスについて、Z値が決定される、（2）Z値の最上位ビットが設定されている場合、 $n/2$ である最大ブロック・サイズのレベルで、 n 個の書き込みデータ要素が交換される、（3）Z値の次の上位ビットが設定されている場合、 $n/4$ である次のブロック・サイズのレベルで、 n 個の書き込みデータ要素が交換される、（4）などとなる。これが、図12において説明される。4つ、8つ、または16個のバンクを含むメモリ・システムの場合、これは、それぞれ2つ、3つ、4つの異なる両方向シャッフル機能のみが必要になり、完全な n 対 n シャッフル・ユニットと比較して著しい単純化であるということの意味する。

【0091】

したがって、やはり簡単に言うと、書き込み動作は、図12のフローチャート1200において要約され得る。まず、アドレス「a」が取得される（1202）。ステップ1204で、パリティおよびZフィールドが導出される。次に、Zフィールドのビット位置（ $\log_2(n) - 1$ ）に基づいて、 $n/2$ のブロック・サイズの粒度で、右側のデータ要素がシャッフルされる（1206）。ステップ1208で、Zフィールドの位置（ $\log_2(n) - 2$ ）に基づいて、 $n/4$ のブロック・サイズの粒度で、再び同じことが発生する、などとなる。このプロセスが繰り返され（破線矢印で示される）、ステップ1210（Zフィールドのビット位置0に基づく、2のブロック・サイズの粒度での、右側のデータ要素のシャッフル）を回答する。

【0092】

新しいシャッフル機能の基本的特性は次のとおりである。

1. それぞれ2つの入力のみをシャッフルする $\log_2(n)$ 個のシャッフル機能（最大で $n * w$ 個の2入力マルチプレクサをそれぞれ必要とし、 w は各データ要素の幅である）

10

20

30

40

50

のみを使用して実装され得る。

2. 各シャッフル機能は、Zベクトルまたはパリティ機能（それぞれ、 $\log_2(n)$ のビット幅である）の1ビットによって制御される。

【0093】

個別の書き込みイネーブル信号が、各書き込みデータ要素に使用される場合、各書き込みデータ要素を含む正しい書き込み信号を適用するために、それらの書き込みイネーブル信号が同じ方法でシャッフルされる必要がある。

【0094】

3. メモリ・バンクから読み取った後のデータ・ベクトル要素の効率的なシャッフル。前述したマッピング方式の「交換」特性のため、書き込みデータに関して上で説明されたシャッフルの概念と同じ概念が、読み取りデータのシャッフルに使用され得る。

10

【0095】

チップ面積またはFPGAのリソースが非常に制限されている場合、 $\log_2(n)$ 個のシャッフルの1セットを、データの読み取りまたは書き込みのいずれかに使用されるメモリ・ポートに使用することもできる。明らかに、そのような使用は、シャッフルされる必要がある読み取りデータ要素または書き込みデータ要素を選択するために、シャッフル段の前に入力マルチプレクサを必要とする。

【0096】

4. パリティ・ビットの各々の計算に使用されるアドレス・ビットの選択。前に示したように、上記の書き込みデータおよび読み取りデータのシャッフルの概念は、パリティ値(X)がオフセット値と完全に同じように順序付けられるということを仮定する。この仮定は、アライメントされたアドレス・シーケンス内の異なる最下位アドレス・ビットが最下位パリティ(X)ビットによってカバーされ、異なる第2のアドレス・ビットが次のパリティ・ビットによってカバーされる、などとなることを必要とする。

20

【0097】

これが当てはまらない例は、図8に示されているマッピング方式のストライド2を伴うアドレス・シーケンス0、2、4、6、8、10、12、および14である。この場合、8つのアドレスの間の $\log_2(n) = \log_2 8 = 3$ つの異なるビット位置は、1、2、および3である。それらのビット位置が、同じ順序でXビット位置に対応しないため、結果として生じるX値は、アドレス0のX=0、アドレス2のX=2、アドレス4のX=4、アドレス6のX=6、アドレス8のX=1、アドレス10のX=3、アドレス12のX=5、およびアドレス14のX=7というように、アドレスと同じ方法で順序付けられない。これらのアドレスは、競合せずに、すべてのメモリ・バンクにわたって完全にインターリーブされるが、上記の簡略化されたシャッフル動作は不十分である。

30

【0098】

この問題は、(1)追加のシャッフル層を追加すること、または(2)パリティ・ベクトルXを生成するために使用されるアドレス・ビットの異なる選択を適用すること、のいずれかによって解決され得る。第1のオプションは、2の累乗のストライドを伴うアドレス・シーケンス内の $\log_2(n)$ 個の異なるビット位置が、Xフィールドのビットが生成される方法を「使用して」アライメントされない、 $(\log_2(n)) - 1$ 個の方法が存在する可能性があるため、 $(\log_2(n)) - 1$ 個の追加の層の追加を必要とする。したがって、図8の例の場合、 $(\log_2(n)) - 1 = 1$ つの追加のシャッフル層が必要になる。

40

【0099】

第2のオプションは、アドレスからXビットが導出される方法を変更する。

【0100】

上記の例では、同じ結果を取得するために、2つのXビット(X0およびX1)の位置が反転され得る。代替として、関連する2の累乗のストライドに関連するアドレス・ビット位置のみから、Xビットが導出され得る。これを実行することによって、前述した書き込みデータおよび読み取りデータのシャッフルの概念と同じ概念が使用され得る。この選

50

択は、それに応じて、パリティが計算されるアドレス・ビットを選択するために、アドレスに適用されるXビット位置ごとにマスク・ベクトルを選択することによって行われ得る（マスクされたアドレス・ビットは、0を与え、XORに基づくパリティ計算を変更しない）。

【0101】

また、このために、アライメントされたアドレス・シーケンス内の異なるビットに関連しない他のすべてのアドレス・ビット位置は、パリティ計算に含まれる必要がない（バンク・オフセットの一部でないすべてのビットがパリティ計算の対象になる必要がある、テーブルなしの方法の場合を除く）。

【0102】

完全性の理由から、図13は、バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするためのメモリ・アクセス・ユニット1400の概略ブロック図を示している。メモリ・アクセス・ユニット1400は、バイナリ・データ・ベクトルの開始アドレスおよびバイナリ・データ・ベクトルの要素の2の累乗のストライドを受信するように適応された受信ユニット1402を備えている。

【0103】

メモリ・アクセス・ユニット1400は、メモリ・バンクの各々について1つのn個のオフセットを決定するように適応された第1の決定モジュール1404をさらに備え、オフセットの各々は、複数のビットレベルのXOR関数を開始アドレスに適用してZベクトルを生成し、マッピング・テーブルにアクセスするためにZベクトルを使用し、バイナリ・データ・ベクトルの2の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることによって決定される。

【0104】

さらに、メモリ・アクセス・ユニット1400は、Zベクトルの2進等価値に応じてn個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定するように適応された第2の決定モジュール1406と、メモリ・ユニットのn個のメモリ・バンク内のバイナリ・データ・ベクトルに並列にアクセスするように適応された基本アクセス・ユニット1408とも備えている。

【0105】

本発明の実施形態は、プログラム・コードを格納することまたは実行することあるいはその両方に適しているプラットフォームに関わらず、事実上、任意の種類のコピュータと一緒に実装されてよい。図14は、一例として、提案された方法に関連するプログラム・コードを実行するのに適しているか、またはメモリ・ユニットに関して本明細書で提案されたアドレス指定方式が実装されてよいか、あるいはその両方であるコンピューティング・システム1500を示している。

【0106】

コンピューティング・システム1500は、適切なコンピュータ・システムの一例にすぎず、コンピュータ・システム1500が上記で示された機能のいずれかを実装されること、または実行すること、あるいはその両方を行うことができるかどうかに関わらず、本明細書に記載された本発明の実施形態の使用または機能の範囲に関してどのような制限も示唆するよう意図されていない。コンピュータ・システム1500には、他の多数の汎用または専用のコンピューティング・システム環境または構成と共に動作できるコンポーネントが存在する。コンピュータ・システム/サーバ1500と共に使用するのに適した周知のコンピューティング・システム、環境、または構成、あるいはその組み合わせの例としては、パーソナル・コンピュータ・システム、サーバ・コンピュータ・システム、シン・クライアント、シック・クライアント、ハンドヘルドまたはラップトップ・デバイス、マイクロプロセッサ・システム、マイクロプロセッサベース・システム、セット・トップ・ボックス、プログラマブル・コンシューマ・エレクトロニクス、ネットワークPC、マイクロコンピュータ・システム、メインフレーム・コンピュータ・システム、およびこれ

10

20

30

40

50

らの任意のシステムまたはデバイスを含む分散クラウド・コンピューティング環境などが挙げられるが、これらに限定されない。コンピュータ・システム/サーバ1500は、コンピュータ・システム1500によって実行されているプログラム・モジュールなどの、コンピュータ・システムによって実行可能な命令との一般的な関連において説明されてよい。通常、プログラム・モジュールは、特定のタスクを実行するか、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、論理、データ構造などを含んでよい。コンピュータ・システム/サーバ1500は、通信ネットワークを介してリンクされたりリモート処理デバイスによってタスクが実行される、分散クラウド・コンピューティング環境で実行されてよい。分散クラウド・コンピューティング環境において、プログラム・モジュールは、メモリ・ストレージ・デバイスを含む、ローカルおよびリモートの両方のコンピュータ・システム・ストレージ媒体に配置されてよい。

10

【0107】

図に示されているように、コンピュータ・システム/サーバ1500は、汎用コンピューティング・デバイスの形態で示されている。コンピュータ・システム/サーバ1500のコンポーネントは、1つまたは複数のプロセッサまたはプロセッシング・ユニット1502、システム・メモリ1504、およびシステム・メモリ1504を含むさまざまなシステム・コンポーネントをプロセッサ1502に結合するバス1506を含むことができるが、これらに限定されない。バス1506は、メモリ・バスまたはメモリ・コントローラ、ペリフェラル・バス、アクセラレーテッド・グラフィックス・ポート、およびさまざまなバス・アーキテクチャのいずれかを使用するプロセッサまたはローカル・バスを含む、複数の種類のバス構造のいずれかのうちの1つまたは複数を表す。例として、そのようなアーキテクチャは、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) ローカル・バス、およびPCI (Peripheral Component Interconnects) バスを含むが、これらに限定されない。コンピュータ・システム/サーバ1500は、通常、さまざまなコンピュータ・システム可読媒体を含む。そのような媒体は、コンピュータ・システム/サーバ1500によってアクセスできる任意の使用可能な媒体であってよく、揮発性および不揮発性媒体、取り外し可能および取り外し不可の媒体を含む。

20

【0108】

システム・メモリ1504は、ランダム・アクセス・メモリ (RAM: random access memory) 1508またはキャッシュ・メモリ1510あるいはその両方などの、揮発性メモリの形態でコンピュータ・システム可読媒体を含んでよい。コンピュータ・システム/サーバ1500は、その他の取り外し可能/取り外し不可、揮発性/不揮発性のコンピュータ・システム・ストレージ媒体をさらに含んでよい。単に例として、取り外し不可、不揮発性の磁気媒体 (図示されておらず、通常は「ハード・ドライブ」と呼ばれる) に対する読み取りと書き込みを行うために、ストレージ・システム1512が提供されてよい。図示されていないが、取り外し可能、不揮発性の磁気ディスク (例えば、「フロッピー (R) ・ディスク」) に対する読み取りと書き込みを行うための磁気ディスク・ドライブ、およびCD-ROM、DVD-ROM、またはその他の光媒体などの取り外し可能、不揮発性の光ディスクに対する読み取りと書き込みを行うための光ディスク・ドライブが提供されてよい。そのような例では、それぞれを、1つまたは複数のデータ媒体インターフェイスによってバス1506に接続することができる。下で詳細に示され、説明されているように、メモリ1504は、本発明の実施形態の機能を実行するように構成された一連の (例えば、少なくとも1つの) プログラム・モジュールを備える少なくとも1つのプログラム製品を含んでよい。

30

40

【0109】

例えば、一連の (少なくとも1つの) プログラム・モジュール1516を含んでいるプログラム/ユーティリティがメモリ1504に格納されてよいが、これに限定されず、オペレーティング・システム、1つまたは複数のアプリケーション・プログラム、その他の

50

プログラム・モジュール、およびプログラム・データも格納されてよい。オペレーティング・システム、1つまたは複数のアプリケーション・プログラム、その他のプログラム・モジュール、およびプログラム・データまたはこれらの組み合わせの各々は、ネットワーク環境の実装を含んでよい。プログラム・モジュール1516は、通常、本明細書に記載された本発明の実施形態の機能または方法あるいはその両方を実行する。

【0110】

また、コンピュータ・システム/サーバ1500は、キーボード、ポインティング・デバイス、ディスプレイ1520などの1つまたは複数の外部デバイス1518、ユーザがコンピュータ・システム/サーバ1500と情報をやりとりできるようにする1つまたは複数のデバイス、またはコンピュータ・システム/サーバ1500が1つまたは複数の他のコンピュータ・システム/サーバ1500と通信できるようにする任意のデバイス（例えば、ネットワーク・カード、モデムなど）、あるいはその組み合わせと通信することもできる。そのような通信は、入出力（I/O：Input/Output）インターフェイス1514を介して行うことができる。さらに、コンピュータ・システム/サーバ1500は、ローカル・エリア・ネットワーク（LAN：local area network）、一般的な広域ネットワーク（WAN：wide area network）、またはパブリック・ネットワーク（例えば、インターネット）、あるいはその組み合わせなどの1つまたは複数のネットワークと、ネットワーク・アダプタ1522を介して通信してよい。図示されているように、ネットワーク・アダプタ1522は、バス1506を介してコンピュータ・システム/サーバ1500の他のコンポーネントと通信してよい。図示されていないが、その他のハードウェア・コンポーネントまたはソフトウェア・コンポーネントあるいはその両方を、コンピュータ・システム/サーバ1500と併用できるということが理解されるべきである。その例として、マイクロコード、デバイス・ドライバ、冗長プロセッシング・ユニット、外部ディスク・ドライブ・アレイ、RAIDシステム、テープ・ドライブ、およびデータ・アーカイブ・ストレージ・システムなどが挙げられるが、これらに限定されない。

【0111】

さらに、バイナリ・データ・ベクトルが一部に格納される複数のn個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするためのメモリ・アクセス・ユニット1400が、バス・システム1506に取り付けられてよい。

【0112】

本発明のさまざまな実施形態の説明は、例示の目的で提示されているが、網羅的であることは意図されておらず、開示された実施形態に制限されない。説明された実施形態の範囲および精神を逸脱することなく多くの変更および変形が可能であることは、当業者にとって明らかである。本明細書で使用された用語は、実施形態の原理、実際の適用、または市場で見られる技術を超える技術的改良を最も適切に説明するため、または他の当業者が本明細書で開示された実施形態を理解できるようにするため選択されている。

【0113】

本発明は、システム、方法、またはコンピュータ・プログラム製品、あるいはその組み合わせとして具現化されてよい。コンピュータ・プログラム製品は、プロセッサに本発明の態様を実行させるためのコンピュータ可読プログラム命令を含んでいるコンピュータ可読ストレージ媒体を含んでよい。

【0114】

この媒体は、電子、磁気、光、電磁気、赤外線、または伝搬媒体用の半導体システムであってよい。コンピュータ可読媒体の例としては、半導体メモリまたは固体メモリ、磁気テープ、取り外し可能フロッピー（R）・ディスク、ランダム・アクセス・メモリ（RAM：random access memory）、読み取り専用メモリ（ROM：read-only memory）、剛体磁気ディスク、および光ディスクが挙げられる。光ディスクの現在の例は、コンパクト・ディスク読み取り専用メモリ（CD-ROM：compact disk-read only memory）、コンパクト・ディスク読み取り/書き込み（CD-R/W：compact disk-read/write）、およびブルーレイディスクを含む。

10

20

30

40

50

【 0 1 1 5 】

コンピュータ可読ストレージ媒体は、命令実行デバイスによって使用するための命令を保持および格納できる有形のデバイスであることができる。コンピュータ可読ストレージ媒体は、例えば、電子ストレージ・デバイス、磁気ストレージ・デバイス、光ストレージ・デバイス、電磁ストレージ・デバイス、半導体ストレージ・デバイス、またはこれらの任意の適切な組み合わせであってよいが、これらに限定されない。コンピュータ可読ストレージ媒体のさらに具体的な例の非網羅的リストは、ポータブル・フロッピー（R）・ディスク、ハード・ディスク、ランダム・アクセス・メモリ（RAM）、読み取り専用メモリ（ROM）、消去可能プログラマブル読み取り専用メモリ（EPROM：erasable programmable read-only memoryまたはフラッシュ・メモリ）、スタティック・ランダム・アクセス・メモリ（SRAM：static random access memory）、ポータブル・コンパクト・ディスク読み取り専用メモリ（CD-ROM：compact disc read-only memory）、デジタル多用途ディスク（DVD：digital versatile disk）、メモリ・スティック、フロッピー（R）・ディスク、パンチカードまたは命令が記録されている溝の中の隆起構造などの機械的にエンコードされるデバイス、およびこれらの任意の適切な組み合わせを含む。本明細書において使用されるとき、コンピュータ可読ストレージ媒体は、それ自体が、電波またはその他の自由に伝搬する電磁波、導波管またはその他の送信媒体を伝搬する電磁波（例えば、光ファイバ・ケーブルを通過する光パルス）、あるいはワイヤを介して送信される電気信号などの一過性の信号であると解釈されるべきではない。

10

【 0 1 1 6 】

本明細書に記載されたコンピュータ可読プログラム命令は、コンピュータ可読ストレージ媒体から各コンピューティング・デバイス/処理デバイスへ、またはネットワーク（例えば、インターネット、ローカル・エリア・ネットワーク、広域ネットワーク、または無線ネットワーク、あるいはその組み合わせ）を介して外部コンピュータまたは外部ストレージ・デバイスへダウンロードされ得る。このネットワークは、銅伝送ケーブル、光伝送ファイバ、無線送信、ルータ、ファイアウォール、スイッチ、ゲートウェイ・コンピュータ、またはエッジ・サーバ、あるいはその組み合わせを備えてよい。各コンピューティング・デバイス/処理デバイス内のネットワーク・アダプタ・カードまたはネットワーク・インターフェイスは、コンピュータ可読プログラム命令をネットワークから受信し、それらのコンピュータ可読プログラム命令を各コンピューティング・デバイス/処理デバイス内のコンピュータ可読ストレージ媒体に格納するために転送する。

20

30

【 0 1 1 7 】

本発明の動作を実行するためのコンピュータ可読プログラム命令は、アセンブラ命令、命令セット・アーキテクチャ（ISA：instruction-set-architecture）命令、マシン命令、マシン依存命令、マイクロコード、ファームウェア命令、状態設定データ、あるいは、Smalltalk、C++などのオブジェクト指向プログラミング言語、および「C」プログラミング言語または同様のプログラミング言語などの従来の手続き型プログラミング言語を含む1つまたは複数のプログラミング言語の任意の組み合わせで記述されたソース・コードまたはオブジェクト・コードであってよい。コンピュータ可読プログラム命令は、ユーザのコンピュータ上で全体的に実行すること、ユーザのコンピュータ上でスタンドアロン・ソフトウェア・パッケージとして部分的に実行すること、ユーザのコンピュータ上およびリモート・コンピュータ上でそれぞれ部分的に実行すること、あるいはリモート・コンピュータ上またはサーバ上で全体的に実行することができる。後者のシナリオでは、リモート・コンピュータは、ローカル・エリア・ネットワーク（LAN）または広域ネットワーク（WAN）を含む任意の種類ネットワークを介してユーザのコンピュータに接続されてよく、または接続は、（例えば、インターネット・サービス・プロバイダを使用してインターネットを介して）外部コンピュータに対して行われてよい。一部の実施形態では、本発明の態様を実行するために、例えばプログラマブル論理回路、フィールドプログラマブル・ゲート・アレイ（FPGA：field-programmable gate arrays）、またはプログラマブル・ロジック・アレイ（PLA：programmable logic arrays）を

40

50

含む電子回路は、コンピュータ可読プログラム命令の状態情報を利用することによって、電子回路をカスタマイズするためのコンピュータ可読プログラム命令を実行してよい。

【0118】

本発明の態様は、本明細書において、本発明の実施形態に従って、方法、装置（システム）、およびコンピュータ・プログラム製品のフローチャート図またはブロック図あるいはその両方を参照して説明される。フローチャート図またはブロック図あるいはその両方の各ブロック、ならびにフローチャート図またはブロック図あるいはその両方に含まれるブロックの組み合わせが、コンピュータ可読プログラム命令によって実装され得るということが理解されるであろう。

【0119】

これらのコンピュータ可読プログラム命令は、コンピュータまたはその他のプログラム可能なデータ処理装置のプロセッサを介して実行される命令が、フローチャートまたはブロック図あるいはその両方のブロックに指定される機能/動作を実施する手段を作り出すべく、汎用コンピュータ、専用コンピュータ、または他のプログラム可能なデータ処理装置のプロセッサに提供されてマシンを作り出すものであってよい。これらのコンピュータ可読プログラム命令は、命令が格納されたコンピュータ可読ストレージ媒体がフローチャートまたはブロック図あるいはその両方のブロックに指定される機能/動作の態様を実施する命令を含んでいる製品を備えるように、コンピュータ可読ストレージ媒体に格納され、コンピュータ、プログラム可能なデータ処理装置、または他のデバイス、あるいはその組み合わせに特定の方式で機能するように指示できるものであってもよい。

【0120】

コンピュータ可読プログラム命令は、コンピュータ、その他のプログラマブル・データ処理装置、あるいは一連の動作可能なステップをコンピュータ上、その他のプログラマブル装置上、またはコンピュータ実装プロセスを生成するその他のデバイス上で実行させる別のデバイスに読み込むこともでき、それによって、コンピュータ上、その他のプログラマブル装置上、または別のデバイス上で実行される命令は、フローチャートまたはブロック図あるいはその両方のブロックにおいて規定された機能/動作を実装する。

【0121】

図内のフローチャートまたはブロック図あるいはその両方は、本発明のさまざまな実施形態に従って、システム、方法、およびコンピュータ・プログラム製品の可能な実装のアーキテクチャ、機能、および動作を示す。これに関連して、フローチャートまたはブロック図内の各ブロックは、規定された論理機能を実装するための1つまたは複数の実行可能な命令を備える、命令のモジュール、セグメント、または部分を表してよい。一部の代替の実装では、ブロックに示された機能は、図に示された順序とは異なる順序で発生してよい。例えば、連続して示された2つのブロックは、実際には、含まれている機能に応じて、実質的に同時に実行されるか、または場合によっては逆の順序で実行されてよい。ブロック図またはフローチャート図あるいはその両方の各ブロック、ならびにブロック図またはフローチャート図あるいはその両方に含まれるブロックの組み合わせは、規定された機能または動作を実行するか、または専用ハードウェアとコンピュータ命令の組み合わせを実行する専用ハードウェアベースのシステムによって実装できるということにも注意する。

【0122】

本明細書で使用される用語は、特定の実施形態を説明することのみを目的としており、本発明を制限することを意図していない。本明細書で使用される単数形「a」、「an」、および「the」は、特に明示的に示されない限り、複数形も含むことが意図されている。「備える」または「備えている」あるいはその両方の用語は、本明細書で使用される場合、記載された機能、整数、ステップ、動作、要素、またはコンポーネント、あるいはその組み合わせの存在を示すが、1つまたは複数のその他の機能、整数、ステップ、動作、要素、コンポーネント、またはこれらのグループ、あるいはその組み合わせの存在または追加を除外していないということが、さらに理解されるであろう。

【0123】

10

20

30

40

50

下の特許請求の範囲内のすべての手段またはステップおよび機能要素の対応する構造、材料、動作、および等価なものは、具体的に請求されるその他の請求された要素と組み合わせる機能を実行するための任意の構造、材料、または動作を含むことが意図されている。本発明の説明は、例示および説明の目的で提示されているが、網羅的であることは意図されておらず、開示された形態での発明に制限されない。本発明の範囲および精神を逸脱することなく多くの変更および変形が可能であることは、当業者にとって明らかである。本発明の原理および実際の適用を最も適切に説明するため、およびその他の当業者が、企図されている特定の用途に適しているようなさまざまな変更を伴う多様な実施形態に関して、本発明を理解できるようにするために、実施形態が選択されて説明された。

【0124】

簡単に言うと、本発明の概念は、バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでいるメモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするための方法として要約されてよく、この方法は、バイナリ・データ・ベクトルの開始アドレスおよびバイナリ・データ・ベクトルの要素の 2 の累乗のストライド s を受信することと、メモリ・バンクの各々について 1 つの n 個のオフセットを決定することと、オフセットの各々が、複数のビットレベルの XOR 関数を開始アドレスに適用して Z ベクトルを生成し、マッピング・テーブルにアクセスするために Z ベクトルを使用し、バイナリ・データ・ベクトルの 2 の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることによって決定される、決定することと、Z ベクトルの 2 進等価値に応じて n 個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定することと、メモリ・ユニットの n 個のメモリ・バンク内のバイナリ・データ・ベクトルに並列にアクセスすることを含む。

【0125】

この方法では、複数のビットレベルの XOR 関数を開始アドレスに適用することは、開始アドレスから、 $\log_2(n)$ 個の部分および $\log_2(n)$ ビット・サイズの Y の部分を選択することと、 $\log_2(n)$ 個の部分の各々についてパリティ・ビットを決定し、 $\log_2(n)$ 次元のパリティ・ベクトル X を生成することと、X と Y の間でビットレベルの XOR 演算を実行して Z ベクトルを生成することを含む。

【0126】

マッピング・テーブルにアクセスするために Z ベクトルを使用することは、マッピング・テーブル内のインデックスとして Z ベクトルを使用することによって、 $\log_2(s)$ 個のエントリを含んでいるマッピング・テーブル内の行を選択することと、マッピング・テーブルのアクセス結果を開始アドレスの部分と結合して n 個のオフセットを取得することを含む。

【0127】

この結合することは、開始アドレスのビット・オフセット $\log_2(n)$ にある $\log_2(n)$ ビットを、決定された $\log_2(n)$ 個のマッピング・テーブル・エントリに置き換え、 $\log_2(n)$ 個のオフセットを生成することを含むとよく、各オフセットは、 $\log_2(n)$ 個のバンクの各 1 つにおいてオフセットとして使用される。

【0128】

n 個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定することは、 $\log_2(n)$ 個の連続するマルチプレクサ段を制御することを含むとよく、マルチプレクサ段の各々は、データ・ベクトルの全幅をカバーするように、複数の 2 入力ビット・マルチプレクサを含む。

【0129】

アクセス動作は、読み取り動作であってよい。

【0130】

アクセス動作は、書き込み動作であってよい。

【0131】

開始アドレスの選択された $\log_2(n)$ 個の部分は、同じサイズに設定されてよい。

10

20

30

40

50

【 0 1 3 2 】

バイナリ・データ・ベクトルは、複数の n 個のメモリ・バンク内の同じサイズの部分に格納されてよい。

【 0 1 3 3 】

開始アドレスの提供、オフセットの決定、ならびにバイナリ・データ・ベクトルの部分のシーケンスの決定、および部分のシーケンスの決定は、 n 個のメモリ・バンクへのバイナリ・データ・ベクトルの書き込み動作の前に実行されてよい。

【 0 1 3 4 】

開始アドレスを提供することと、オフセットを決定することと、バイナリ・データ・ベクトルの部分のシーケンスを決定することとは、 n 個のメモリ・バンクに対するバイナリ・データ・ベクトルの読み取り動作の前に実行されてよく、部分のシーケンスを決定することは、 n 個のメモリ・バンクに対するバイナリ・データ・ベクトルの読み取り動作の後に実行されてよい。

10

【 0 1 3 5 】

メモリ・ユニット内のバイナリ・データ・ベクトルにアクセスするためのメモリ・アクセス・ユニットが、バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでおり、メモリ・アクセス・ユニットは、バイナリ・データ・ベクトルの開始アドレスおよびバイナリ・データ・ベクトルの要素の 2 の累乗のストライド s を受信するように適応された受信ユニットと、メモリ・バンクの各々について 1 つの n 個のオフセットを決定するように適応された第 1 の決定モジュールであって、オフセットの各々が、複数のビットレベルの XOR 関数を開始アドレスに適用して Z ベクトルを生成し、マッピング・テーブルにアクセスするために Z ベクトルを使用し、バイナリ・データ・ベクトルの 2 の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることによって決定される、第 1 の決定モジュールと、Z ベクトルの 2 進等価値に応じて n 個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定するように適応された第 2 の決定モジュールと、メモリ・ユニットの n 個のメモリ・バンク内のバイナリ・データ・ベクトルに並列にアクセスするように適応された基本アクセス・ユニットとを備える。

20

【 0 1 3 6 】

第 1 の決定モジュールによって複数のビットレベルの XOR 関数を開始アドレスに適用することは、開始アドレスから、 $\log_2(n)$ 個の部分および $\log_2(n)$ ビット・サイズの Y の部分を選択することと、 $\log_2(n)$ 個の部分の各々についてパリティ・ビットを決定し、 $\log_2(n)$ 次元のパリティ・ベクトル X を生成することと、X と Y の間でビットレベルの XOR 演算を実行して Z ベクトルを生成することとを含んでよい。

30

【 0 1 3 7 】

第 1 の決定モジュールによってマッピング・テーブルにアクセスするために Z ベクトルを使用することは、マッピング・テーブル内のインデックスとして Z ベクトルを使用することによって、 $\log_2(s)$ 個のエントリを含んでいるマッピング・テーブル内の行を選択することと、マッピング・テーブルのアクセス結果を開始アドレスの部分と結合して n 個のオフセットを取得することとを含んでよい。

40

【 0 1 3 8 】

第 1 の決定モジュールの結合は、開始アドレスのビット・オフセット $\log_2(n)$ にある $\log_2(n)$ ビットを、決定された $\log_2(n)$ 個のマッピング・テーブル・エントリに置き換え、 $\log_2(n)$ 個のオフセットを生成することを含んでもよく、各オフセットは、 $\log_2(n)$ 個のバンクの各 1 つにおいてオフセットとして使用される。

【 0 1 3 9 】

第 2 の決定モジュールは、 n 個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定するときに、 $\log_2(n)$ 個の連続するマルチプレクサ段を制御するように適応されてもよく、マルチプレクサ段の各々は、データ・ベクトルの全幅をカバーするように、複数の 2 入力ビット・マルチプレクサを含む。

50

【0140】

アクセス動作は、読み取り動作であってよい。

【0141】

アクセス動作は、書き込み動作であってよい。

【0142】

開始アドレスの選択された $\log_2(n)$ 個の部分は、同じサイズに設定されてよい。

【0143】

バイナリ・データ・ベクトルは、複数の n 個のメモリ・バンク内の同じサイズの部分に格納されてよい。

【0144】

開始アドレスの提供、オフセットの決定、ならびにバイナリ・データ・ベクトルの部分のシーケンスの決定、および部分のシーケンスの決定は、 n 個のメモリ・バンクへのバイナリ・データ・ベクトルの書き込み動作の前に実行されてよい。

【0145】

開始アドレスを提供することと、オフセットを決定することと、バイナリ・データ・ベクトルの部分のシーケンスを決定することとは、 n 個のメモリ・バンクに対するバイナリ・データ・ベクトルの読み取り動作の前に実行されてよく、部分のシーケンスを決定することは、 n 個のメモリ・バンクに対するバイナリ・データ・ベクトルの読み取り動作の後に実行されてよい。

【0146】

バイナリ・データ・ベクトルが一部に格納される複数の n 個のメモリ・バンクを含んでいるメモリ・コンポーネント内のバイナリ・データ・ベクトルにアクセスするコンピュータ・プログラム製品であって、コンピュータ・プログラム製品が、プログラム命令が具現化されているコンピュータ可読ストレージ媒体を備えており、プログラム命令が、1つまたは複数のコンピューティング・システムに、バイナリ・データ・ベクトルの開始アドレスおよびバイナリ・データ・ベクトルの要素の2の累乗のストライド s を受信することと、メモリ・バンクの各々について1つの n 個のオフセットを決定することであって、オフセットの各々が、複数のビットレベルの XOR 関数を開始アドレスに適用して Z ベクトルを生成することと、マッピング・テーブルにアクセスするために Z ベクトルを使用することと、バイナリ・データ・ベクトルの2の累乗のストライドに従ってマッピング・テーブルのアクセス結果をシフトすることとによって決定される、決定することと、 Z ベクトルの2進等価値に応じて n 個のメモリ・バンク内のバイナリ・データ・ベクトルの部分のシーケンスを決定することと、メモリ・ユニットの n 個のメモリ・バンク内のバイナリ・データ・ベクトルに並列にアクセスすることとを実行させるように、1つまたは複数のコンピューティング・システムまたはコントローラによって実行可能である、コンピュータ・プログラム製品。

10

20

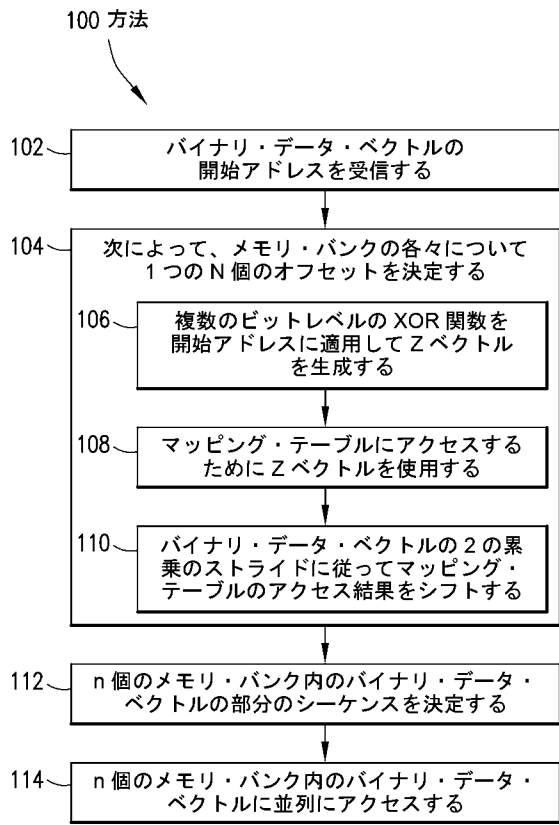
30

40

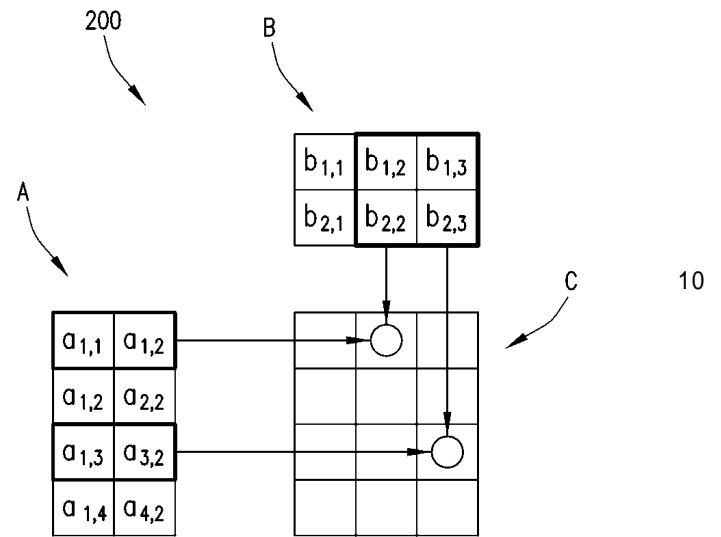
50

【図面】

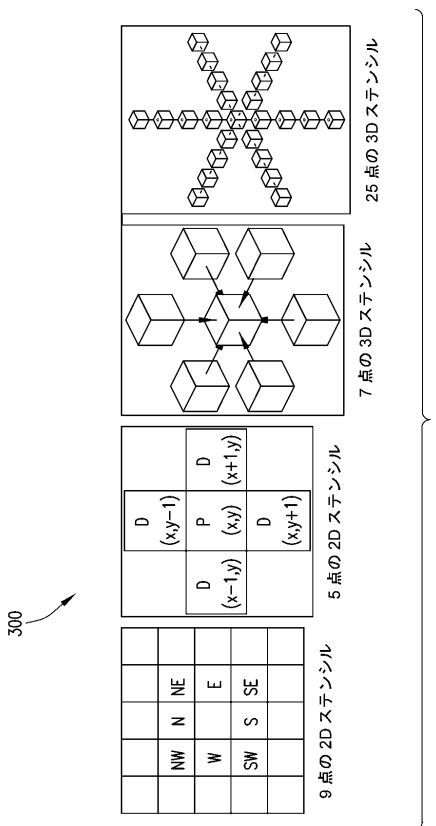
【図 1】



【図 2】



【図 3】



【図 4】

400

27	27	91	155	219	283	347	411	475
26	90	26	218	154	346	282	474	410
25	153	217	25	89	409	473	281	345
24	216	152	88	24	472	408	344	280
23	343	279	471	407	87	23	215	151
22	278	342	406	470	22	86	150	214
21	469	405	341	277	213	149	85	21
20	404	468	276	340	148	212	20	84
19	83	19	211	147	339	275	467	403
18	18	82	146	210	274	338	402	466
17	209	145	81	17	465	401	337	273
16	144	208	16	80	400	464	272	336
15	399	463	271	335	143	207	15	79
14	462	398	334	270	206	142	78	14
13	269	333	397	461	13	77	141	205
12	332	268	460	396	76	12	204	140
11	139	203	11	75	395	459	267	331
10	202	138	74	10	458	394	330	266
9	9	73	137	201	265	329	393	457
8	72	8	200	136	328	264	456	392
7	455	391	327	263	199	135	71	7
6	390	454	262	326	134	198	6	70
5	325	261	453	389	69	5	197	133
4	260	324	388	452	4	68	132	196
3	195	131	67	3	451	387	323	259
2	130	194	2	66	386	450	258	322
1	65	1	193	129	321	257	449	385
0	0	64	128	192	256	320	384	448

オフセット バンク0 バンク1 バンク2 バンク3 バンク4 バンク5 バンク6 バンク7

10

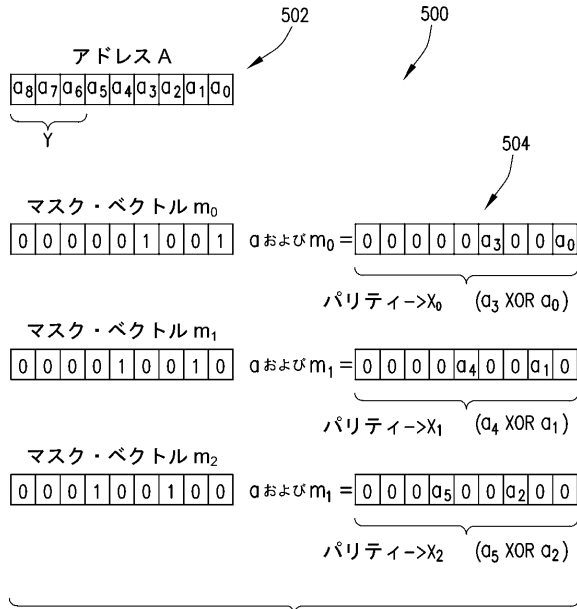
20

30

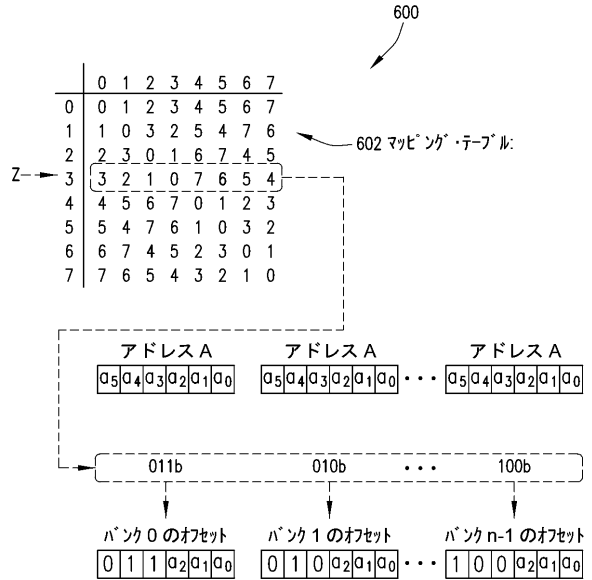
40

50

【図5】

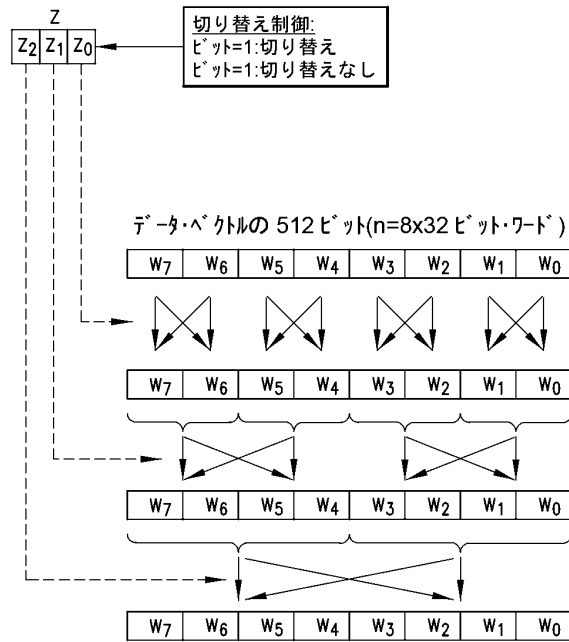


【図6】

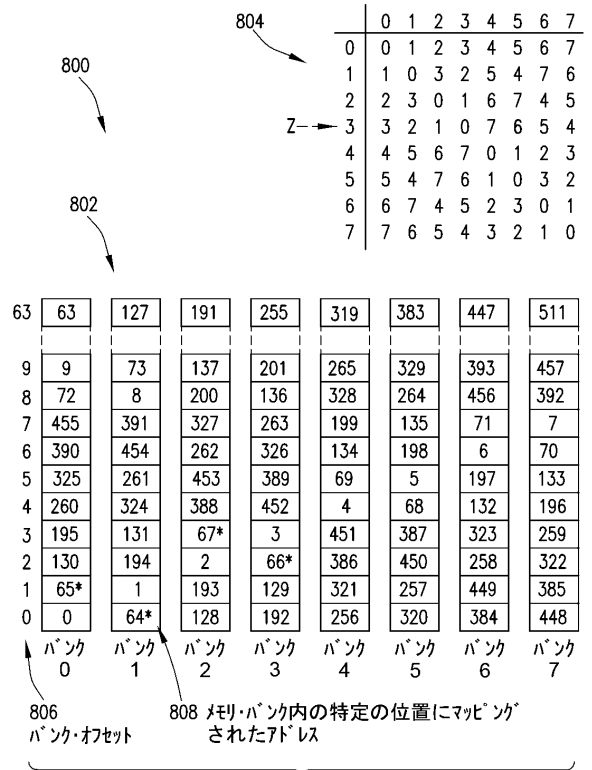


10

【図7】



【図8】



20

30

40

【図 9】

900

63	63	127	191	255	319	383	447	511
9	9	73	137	201	265	329	393	457
8	72	8	200	136	328	264	456	392
7	455	391	327	263	199	135	71	7
6	390	454	262	326	134	198	6	70
5	325	261	453	389	69	5	197	133
4	260	324	388	452	4	68	132	196
3	195	131	67	3	451	387	323	259
2	130	194	2	66	386	450	258	322
1	65	1	193	129	321	257	449	385
0	0	64	128	192	256	320	384	448

バンク0 バンク1 バンク2 バンク3 バンク4 バンク5 バンク6 バンク7

【図 10】

1000

63	63	127	191	255	319	383	447	511
59	315	379	443	507	59	123	187	251
51	371	307	499	435	115	51	243	179
43	427	491	299	363	171	235	43	107
35	483	419	355	291	227	163	99	35
27	27	91	155	219	283	347	411	475
19	83	19	211	147	339	275	467	403
11	139	203	11	75	395	459	267	331
3	195	131	67	3	451	387	323	259
0	0	64	128	192	256	320	384	448

バンク0 バンク1 バンク2 バンク3 バンク4 バンク5 バンク6 バンク7

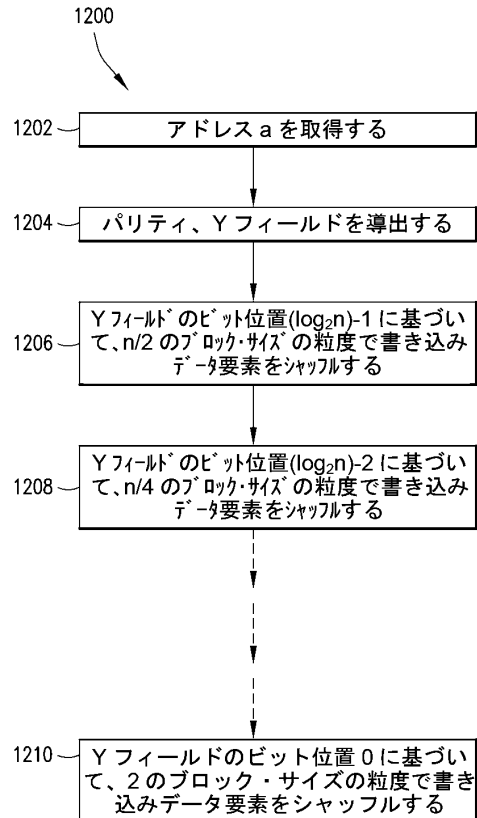
10

20

【図 11】



【図 12】

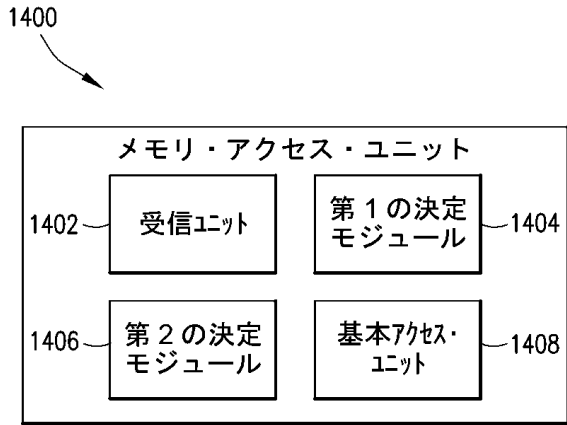


30

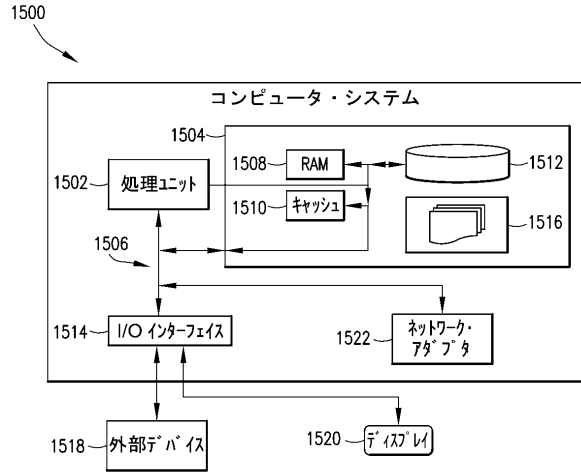
40

50

【図 13】



【図 14】



10

20

30

40

50

フロントページの続き

- (72)発明者 ファン ルンテレン、ヤン
スイス 8 8 0 3 リュシュリコン ゾイマーシュトラーク 4
審査官 漆原 孝治
- (56)参考文献 特開平 0 2 - 1 2 7 7 6 8 (J P , A)
特開 2 0 0 2 - 0 7 3 4 1 2 (J P , A)
特開 2 0 0 0 - 1 6 3 3 1 6 (J P , A)
- (58)調査した分野 (Int.Cl. , D B 名)
G 0 6 F 1 7 / 1 6