

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5298393号
(P5298393)

(45) 発行日 平成25年9月25日(2013.9.25)

(24) 登録日 平成25年6月28日(2013.6.28)

(51) Int.Cl. F I
G 0 6 F 3 / 0 6 (2006.01) G 0 6 F 3 / 0 6 3 0 5 C

請求項の数 26 (全 19 頁)

<p>(21) 出願番号 特願2010-534958 (P2010-534958) (86) (22) 出願日 平成20年11月20日 (2008.11.20) (65) 公表番号 特表2011-504269 (P2011-504269A) (43) 公表日 平成23年2月3日 (2011.2.3) (86) 国際出願番号 PCT/US2008/012944 (87) 国際公開番号 W02009/070235 (87) 国際公開日 平成21年6月4日 (2009.6.4) 審査請求日 平成23年11月17日 (2011.11.17) (31) 優先権主張番号 60/989,670 (32) 優先日 平成19年11月21日 (2007.11.21) (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 502188642 マーベル ワールド トレード リミテッ ド バルバドス国 ビービー14027, セン トマイケル、ブリトンズ ヒル、ガンサイ トロード、エル ホライズン (74) 代理人 110000877 龍華国際特許業務法人 (72) 発明者 ブルーチ、アービンド アメリカ合衆国、95054 カリフォル ニア州、サンタ クララ、マーベル レー ン 5488 マーベル セミコンダクタ ー インコーポレイテッド内 審査官 坂東 博司 最終頁に続く</p>
---	--

(54) 【発明の名称】 並列リードソロモンRAID (RS-RAID) アーキテクチャ、デバイス、および方法

(57) 【特許請求の範囲】

【請求項1】

分散データ記憶デバイスであって、
 複数の並列データ記憶クラスタであって、
 複数のデータシンボルを記憶する複数のデータ記憶デバイス、および、
 前記複数のデータ記憶デバイスのそれぞれに結合し、前記複数のデータシンボルからローカルな中間チェックサムを計算し、前記ローカルな中間チェックサムを含む複数の中間チェックサムに基づいて前記複数のデータシンボルのチェックサムを計算するよう構成されたRAIDコントローラを含む、複数の並列データ記憶クラスタと、
 前記複数の並列データ記憶クラスタのそれぞれに結合し、前記複数の中間チェックサムを前記複数の並列データ記憶クラスタに配信する通信ファブリックとを備える分散データ記憶デバイス。

【請求項2】

前記RAIDコントローラは、前記複数のデータシンボルの重み付き和に基づいて前記ローカルな中間チェックサムを計算する、請求項1に記載の分散データ記憶デバイス。

【請求項3】

前記複数の並列データ記憶クラスタの複数のRAIDコントローラの少なくとも1つは、他の並列データ記憶クラスタから前記中間チェックサムを受信し、前記複数の並列なデータ記憶クラスタの複数のデータ記憶デバイスのうちの1つに記憶するためのチェックサムを、前記受信した中間チェックサムと、前記ローカルな中間チェックサムとに基づいて

10

20

生成する、請求項 1 又は 2 に記載の分散データ記憶デバイス。

【請求項 4】

前記 R A I D コントローラは、前記チェックサムをデータ記憶デバイスに記憶するようさらに構成される請求項 1 から 3 のいずれか 1 項に記載の分散データ記憶デバイス。

【請求項 5】

前記通信ファブリックは、前記複数のデータシンボルを含む複数の情報シンボルを前記複数の並列データ記憶クラスタに配信するようさらに構成される請求項 1 から 4 のいずれか 1 項に記載の分散データ記憶デバイス。

【請求項 6】

前記 R A I D コントローラは、前記複数の情報シンボルのサブセットから前記複数のデータシンボルを選択する R A I D コントロールユニットをさらに備える請求項 5 に記載の分散データ記憶デバイス。

10

【請求項 7】

前記 R A I D コントロールユニットは、前記複数のデータ記憶デバイスから作動可能なデータ記憶デバイスのリストを確定する記憶デバイス故障センスユニットをさらに備える請求項 6 に記載の分散データ記憶デバイス。

【請求項 8】

前記 R A I D コントロールユニットは、作動可能なデータ記憶デバイスの前記リストに基づいてデータ回復行列を計算するようさらに構成される請求項 7 に記載の分散データ記憶デバイス。

20

【請求項 9】

前記通信ファブリックは、作動可能なデータ記憶デバイスの前記リストが変化すると、前記データ回復行列を前記複数の並列データ記憶クラスタに配信するようさらに構成される請求項 8 に記載の分散データ記憶デバイス。

【請求項 10】

前記 R A I D コントローラは、作動可能なデータ記憶デバイスのリストに基づいて前記中間チェックサムを計算するよう構成された中間和デバイスをさらに備える請求項 6 に記載の分散データ記憶デバイス。

【請求項 11】

前記中間和デバイスは、データシンボルが変化すると、前記ローカルな中間チェックサムを更新するよう構成された再計算器をさらに備える請求項 10 に記載の分散データ記憶デバイス。

30

【請求項 12】

前記 R A I D コントロールユニットは、前記データ回復行列を計算するために、添加された単位行列とファンデルモンデ行列の逆行列を求めるようさらに構成される請求項 8 に記載の分散データ記憶デバイス。

【請求項 13】

中間和デバイスは、前記データ回復行列に基づく中間データシンボル、ならびに、前記記憶デバイスから読出されるリードデータシンボルおよびリードチェックサムシンボルの少なくとも一方のベクトルを計算するようさらに構成される請求項 12 に記載の分散データ記憶デバイス。

40

【請求項 14】

前記 R A I D コントロールユニットは、前記添加された単位行列とファンデルモンデ行列の行の集合ならびに作動可能なデータ記憶デバイスの前記リストに基づくベクトルを選択し、前記データ回復行列を形成するために前記行の集合の逆行列を求めるようさらに構成される請求項 12 に記載の分散データ記憶デバイス。

【請求項 15】

前記 R A I D コントローラは、中間和デバイスに中間データシンボルを計算させるメッセージを前記複数の並列データ記憶クラスタに送出するようさらに構成される請求項 14 に記載の分散データ記憶デバイス。

50

【請求項 1 6】

前記通信ファブリックは、前記中間データシンボルを前記複数の並列データ記憶クラスタのそれぞれに配信するようさらに構成される請求項 1 5 に記載の分散データ記憶デバイス。

【請求項 1 7】

前記 R A I D コントローラは、複数の中間データシンボルから、回復されるデータシンボルを計算するようさらに構成される請求項 1 6 に記載の分散データ記憶デバイス。

【請求項 1 8】

前記ローカルな中間チェックサムは、前記 R A I D コントローラに割当てられた前記複数のデータ記憶デバイスに対応する前記複数のデータシンボルの部分集合の部分 and (partial sum) である請求項 1 から 1 7 のいずれか 1 項 に記載の分散データ記憶デバイス。

10

【請求項 1 9】

前記中間データシンボルは、前記 R A I D コントローラに割当てられた前記データ記憶デバイスに対応する前記データシンボルの部分集合および前記リードチェックサムシンボルの部分 and である請求項 1 3 に記載の分散データ記憶デバイス。

【請求項 2 0】

データ記憶方法であって、

複数のデータ記憶デバイスを複数の並列データ記憶クラスタに割当てること、

前記複数の並列データ記憶クラスタのそれぞれに複数のデータシンボルを記憶すること

20

、
前記複数の並列データ記憶クラスタのそれぞれに含まれる複数の R A I D コントローラのそれぞれが、前記複数のデータシンボルの重み付き和から中間チェックサムを計算すること、

前記複数の R A I D コントローラの少なくとも 1 つが、前記複数の R A I D コントローラの他の R A I D コントローラのそれぞれから前記中間チェックサムを受信すること、

前記複数の R A I D コントローラの少なくとも 1 つが、計算した前記中間チェックサムおよび前記他の R A I D コントローラのそれぞれから受信した前記中間チェックサムに基づいて前記複数のデータシンボルのチェックサムを計算すること、および、

前記複数のデータ記憶デバイスの少なくとも 1 つに前記チェックサムを記憶することを含むデータ記憶方法。

30

【請求項 2 1】

情報シンボルの部分集合から前記複数のデータシンボルを選択すること、および、

通信ファブリックを使用して、前記複数の並列データ記憶クラスタに前記選択されたデータシンボルを配信することをさらに含む請求項 2 0 に記載のデータ記憶方法。

【請求項 2 2】

前記複数の並列データ記憶クラスタのそれぞれの中の作動可能なデータ記憶デバイスの集合を検知することをさらに含む請求項 2 1 に記載のデータ記憶方法。

【請求項 2 3】

作動可能なデータ記憶デバイスの前記集合に基づくデータ回復行列を、前記複数の並列データ記憶クラスタのそれぞれに配信することをさらに含む請求項 2 2 に記載のデータ記憶方法。

40

【請求項 2 4】

前記データ回復行列に基づく中間データシンボルならびに前記記憶デバイスから読出されるリードデータシンボルおよびリードチェックサムシンボルの少なくとも一方のベクトルを計算することをさらに含む請求項 2 3 に記載のデータ記憶方法。

【請求項 2 5】

前記中間データシンボルを前記複数の並列データ記憶クラスタのそれぞれに配信すること、および、

前記中間データシンボルの和から、回復されるデータシンボルを計算することをさらに含む請求項 2 4 に記載のデータ記憶方法。

50

【請求項 26】

分散データアーキテクチャにおける誤り訂正のための方法であって、
複数のデータ記憶デバイスを複数の並列データ記憶クラスタに割当てる構成行列を读出すこと、

複数のデータシンボルを前記複数の並列データ記憶クラスタに記憶すること、

前記複数の並列データ記憶クラスタのそれぞれに含まれる複数の R A I D コントローラのそれぞれが、前記割当てられたデータ記憶デバイスに記憶された前記複数のデータシンボルから中間チェックサムを計算することであって、それにより、少なくとも1つの中間チェックサムがそれぞれのデータ記憶クラスタについて計算される、計算すること、

前記複数の R A I D コントローラの少なくとも1つが、前記複数の R A I D コントローラの他の R A I D コントローラのそれぞれから前記中間チェックサムを受信すること、

チェックサムを形成するために、前記複数の R A I D コントローラの少なくとも1つが、計算した前記中間チェックサムおよび前記他の R A I D コントローラのそれぞれから受信した前記中間チェックサムを合計すること、および、

前記チェックサムを少なくとも1つのデータ記憶デバイスに記憶する方法。

【発明の詳細な説明】

【技術分野】

【0001】

本出願は、2007年11月21日に出願された米国仮出願第60/989,670号「Parallel RAID Implementation for RAID6 and Reed-Solomon Code」の利益を主張し、参照によりその全体が本明細書に組み込まれる全ての引用された参考文献を含む。

【背景技術】

【0002】

RAID (redundant array of inexpensive disks) アーキテクチャは、ハードディスクなどのデータ記憶ユニットのグループを使用して、耐故障性を有するデータ記憶装置を提供する。RAID アーキテクチャは、誤りおよびディスク故障から情報を保護するために、前方誤り訂正 (forward error correction) (FEC) コードおよび予備のデータ記憶ユニットを使用する。情報シンボルは、ビット、バイト、またはワードであってよい。情報シンボルは、符合化されて、データおよびチェックサムまたはパリティシンボルを含むコードシンボルを形成しうる。組織的な前方誤り訂正コードの場合、情報シンボルは、コードシンボルのデータシンボル部分において明示的に表されうる。

【0003】

リードソロモンコードは、チェックサムシンボルの数に等しい記憶ユニットの故障の数を許容するために、RAID アーキテクチャ (RS-RAID) で使用されうる。たとえば、データ用に20の記憶ユニットを、チェックサム用に4つの記憶ユニットを割当てる4重誤り訂正 RS-RAID アーキテクチャは、4つを含む4つまでの記憶デバイスにおける故障を許容しうる。

【0004】

RS-RAID アーキテクチャは、通常、データ記憶ユニットに書込まれるデータシンボルを保護するために単一の RAID コントローラを使用する。単一の RAID コントローラが使用されて、チェックサム、符号化、および復号化計算を実施するとき、RAID アーキテクチャのスループットまたはデータ記憶量および取出し速度は、RAID でなくかつ耐故障性を有さないデータ記憶アーキテクチャと比べて低減される可能性がある。

【発明の概要】

【発明が解決しようとする課題】

【0005】

したがって、高スループットで耐故障性を有する分散型のデータ記憶アーキテクチャが望ましい場合がある。

【課題を解決するための手段】

【0006】

高性能記憶アーキテクチャでは、複数のRAIDコントローラは、通信ファブリックと呼ぶ通信経路の共通集合を通じて互いに通信してもよい。通信ファブリックは、RAIDコントローラと所与のRAIDコントローラに割当てられた記憶デバイスとの間の通信経路と比較して高いレーテンシを有する可能性がある。レーテンシの高い通信ファブリックは、RAIDコントローラ間のデータ、メッセージ、構成などのトラフィックが、耐故障性を有する分散データ記憶装置のタスクに整合しなければ、RAIDデータ記憶アーキテクチャのスループットを減少させる可能性がある。通信ファブリックと割当てられたデータ記憶デバイスの集合との間に介在してもよい、それぞれのRAIDコントローラは、データ記憶アーキテクチャのノードと呼ばれてもよい。RAIDコントローラおよび割当てられたデータ記憶デバイスは、データ記憶クラスタと呼ばれてもよい。

10

【0007】

リードソロモン(Reed-Solomon)RAID(RS-RAID)アーキテクチャは、冗長なデータ記憶デバイスを含むことによって、ハードディスクなどの記憶デバイスに書込まれ、また、記憶デバイスから読出される情報シンボルを保護しうる。mのチェックサムデバイスを使用するRS-RAIDアーキテクチャは、データ記憶デバイスのm程度の同時故障を許容しうる。mのチェックサムシンボルは c_1, c_2, \dots, c_m で示されてもよい。RS-RAIDアーキテクチャはまた、 d_1, d_2, \dots, d_n で示す情報保持またはデータシンボル用の数nのデータ記憶デバイスを含みうる。

20

【0008】

チェックサムおよびデータ記憶デバイスは、データおよびチェックサムシンボルを、ビット、バイト、ワードなどとして記憶してもよい。リードソロモン(RS)コードなどのあるタイプの前方誤り訂正コード(forward error correction codes)(FEC)は、通常、バイトを使用することが留意されてもよい。たとえば、RSコードは、255のバイトブロック内の233のデータバイトと32のチェックサムバイトに233の情報バイトを符合化するブロックなどのバイトブロックに作用しうる。

【0009】

RS-RAIDアーキテクチャは、対応するデータ記憶デバイス D_1, D_2, \dots, D_n によって保持されるデータシンボル d_1, d_2, \dots, d_n を使用して、i番目のチェックサムデバイス C_i に記憶されるチェックサムシンボル c_i を計算しうる。RS-RAIDアーキテクチャは、 $D_1, D_2, \dots, D_n, C_1, C_2, \dots, C_m$ の記憶デバイスのうちの任意のm以下の記憶デバイスが故障する場合、故障したデバイスのうちの任意のデバイスのコンテンツが、損なわれていないまたは故障していないデバイスから再構築されうるように各 c_i ($1 \leq i \leq m$)を確定しうる。RS-RAIDアーキテクチャは、ファンデルモンデ行列の特性により耐故障性を有する演算を提供することができ、ファンデルモンデ行列は、チェックサムシンボルを計算し維持し、記憶デバイスから読出したデータおよびチェックサムシンボルから情報を回復するのに使用される。RS-RAIDコントローラは、記憶デバイスが故障しても、添加されるかまたは拡大されたファンデルモンデ行列と単位行列の($n \times n$)部分の逆行列を計算することによって、記憶デバイスにおけるデータおよび/またはチェックサムシンボルを回復しうる。

30

40

【0010】

チェックサムシンボルを生成するために、RS-RAIDアーキテクチャは、データシンボルを、ファンデルモンデ行列の要素で重み付けし、式1によって線形関数 F_i を使用して重み付けされたデータシンボルを合計しうる。関数 F_i は、ファンデルモンデ行列の要素のi番目の行から得られうるため、 $F_i = [f_{i,1}; f_{i,2}; \dots; f_{i,n}]^T$ である。

【数 1】

$$c_i = \sum_{j=1}^n d_j f_{i,j} \quad \text{式 1}$$

換言すれば、データおよびチェックサムシンボルが、それぞれ、 $(n \times 1)$ 次元および $(m \times 1)$ 次元ベクトル $D = [d_1, d_2, \dots, d_n]^T$ および $C = [c_1, c_2, \dots, c_m]^T$ として表され、また、関数 F_i が行列 F の行として表される場合、RS-R A I Dアーキテクチャは、チェックサムシンボルを式 2 a によって符合化しうる。

$$C = F D \quad (\text{式 2 a})$$

式 2 a は、

10

【数 2】

$$\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} = \begin{bmatrix} f_{1,1} & f_{1,2} & \dots & f_{1,n} \\ f_{2,1} & f_{2,2} & \dots & f_{2,n} \\ \vdots & \vdots & & \vdots \\ f_{m,1} & f_{m,2} & \dots & f_{m,n} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix} \quad \text{式 2 b}$$

に等しい。

【0 0 1 1】

20

有利に設計されたRS-R A I D F E Cコードの場合、 F 行列は、要素： $f_{i,j} = j^{i-1}$ を有する $(m \times n)$ ファンデルモンデ行列でありうる。式中、インデックス $i = 1, 2, \dots, m$ および $j = 1, 2, \dots, n$ は、それぞれ、ファンデルモンデ行列の行および列に対応し、代数演算は、ガロア体の特性を使用して実施される。たとえば、 (3×4) ファンデルモンデ行列は、

【数 3】

$$F = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 9 & 16 \end{bmatrix} \quad \text{式 3}$$

30

として書かれうる。

【0 0 1 2】

誤りを含む可能性があるコードワードまたはコードシンボルから $(n \times 1)$ 情報ベクトル

【数 4】

$$\tilde{D} = [\tilde{d}_1, \tilde{d}_2, \dots, \tilde{d}_n]^T$$

を回復するために、並列RS-R A I Dアーキテクチャは、ファンデルモンデ行列および I で示す添加された $((n-m) \times (n-m))$ 単位行列を含む拡大されるかまたは分割された行列 A の逆行列を求め、 A の逆行列に、作動可能な記憶ユニットから読出されるデータおよびチェックサムシンボルの集合 D および C をそれぞれ右から乗算し($post-multiply$)うる。記号的には、回復される情報ベクトル

40

【数 5】

$$\tilde{D}$$

は、

【数 6】

$$\tilde{D} = \text{Inv}(A) \cdot [E]$$

から得られうる。式中、拡大行列は

【数 7】

$$A \Delta \begin{bmatrix} I \\ \dots \\ F \end{bmatrix}$$

であり、

【数 8】

$$E \Delta \begin{bmatrix} D \\ \dots \\ C \end{bmatrix}$$

は、拡大されたデータおよびチェックサムシンボルベクトルである。表記 $\text{Inv}(A)$ は、正則な $(n \times n)$ 正方行列を形成する A の行の部分集合の逆行列などの A に基づく逆行列をもたらし、また、以下で述べるように、

【数 9】

$$[E]$$

で示す列行列 E の n の行の、対応する選択されるかまたは選抜された集合に共形的である関数であると理解されてもよい。 A 行列の逆行列を求めるプロセスは、 A の行の選択された集合の反転とみなされてもよく、選択は、作動可能なデータ記憶デバイスのリストおよび行列にベクトルを掛ける計算における共形性についての要件によって確定される。 $(n + m) \times n$ 拡大行列 A の n の行の全ての部分集合は、 F がファンデルモンデ行列であるため反転可能であることが留意されてもよい。

【0013】

拡張形態では、式

【数 10】

$$\tilde{D} = \text{Inv}(A) \cdot [E]$$

は、

【数 11】

$$\begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \vdots \\ \tilde{d}_n \end{bmatrix} = \text{Inv} \left(\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & 1 \\ \hline 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 2^{m-1} & 3^{m-1} & \dots & n^{m-1} \end{bmatrix} \right) \cdot \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \\ \hline c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} \quad \text{式 4}$$

として表されうる。式中、共形性は、行列 A の選択された部分を反転する前に、 E および A の対応する行を選択することによって実施される。

【0014】

換言すれば、 $RS - RAID$ アーキテクチャの各記憶デバイスは、拡大行列 A の行および列ベクトル $E = [d_1, d_2, \dots, d_n, c_1, c_2, \dots, c_m]^T$ の対応する要素によって表されうる。 m の冗長記憶デバイス中のいずれの冗長記憶デバイスも故障しない場合、回復される情報シンボルは、 A の n の行の任意の部分集合および E の n の対応する要素を選択することによって確定されて、データ回復行列として記述されてもよい正方行列

10

20

30

40

50

A'および対応するデータ記憶ユニットから読出されるデータのベクトル

【数 1 2】

$$E' = [E].$$

が形成されうる。換言すれば、 $\text{Inv}(A) = (A')^{-1}$ であり、また、

【数 1 3】

$$\tilde{D} = (A')^{-1} \cdot E'$$

である。たとえば、4 + 2 R S - R A I Dアーキテクチャの場合、回復されるかまたは復号化されるデータ

10

【数 1 4】

$$\tilde{D}$$

は、

【数 1 5】

$$\begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \tilde{d}_3 \\ \tilde{d}_4 \end{bmatrix} = \text{Inv} \begin{pmatrix} \begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 2 & 3 & 4 \end{matrix} \end{pmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ \hline e_1 \\ \hline e_2 \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{bmatrix} \quad \text{式 5}$$

20

によって、拡大されたファンデルモンデ行列の最初の4つの行および記憶デバイスアレ
イから読出されたデータおよびチェックサム
の最初の4つのエンティティから抽出された
回復されるデータシンボルのベクトルでありうる。

【0 0 1 5】

たとえば、第3の、第5の、または、第3と第5の両方のデータ記憶デバイスが故障する
場合、

30

【数 1 6】

$$\tilde{D},$$

は、以下の通りに、作動可能なデバイスに対応する4つの行を選択することによって、
E'から回復されうる。

【数 1 7】

$$\begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \tilde{d}_3 \\ \tilde{d}_4 \end{bmatrix} = \text{Inv} \begin{pmatrix} \begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ \hline 1 & 2 & 3 & 4 \end{matrix} \end{pmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \\ \hline \cancel{d_3} \\ d_4 \\ \hline c_1 \\ \hline e_2 \end{bmatrix} \quad \text{式 6 a}$$

40

$$= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} d_1 \\ d_2 \\ d_4 \\ c_1 \end{bmatrix}$$

【数 1 8】

$$\begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \tilde{d}_3 \\ \tilde{d}_4 \end{bmatrix} = Inv \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ \hline \hline \hline \hline \\ 1 & 2 & 3 & 4 \end{pmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad \text{式 6 b}$$

$$= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{bmatrix} \quad 10$$

【数 1 9】

$$\begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \tilde{d}_3 \\ \tilde{d}_4 \end{bmatrix} = Inv \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \hline \hline \hline \hline \\ 0 & 0 & 0 & 1 \\ \hline \hline \hline \hline \\ 1 & 2 & 3 & 4 \end{pmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \\ \varepsilon_3 \\ d_4 \\ \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad \text{式 6 c}$$

$$= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}^{-1} \begin{bmatrix} d_1 \\ d_2 \\ d_4 \\ \varepsilon_2 \end{bmatrix} \quad 20$$

式中、2重末梢線は記憶デバイスの故障を示し、1重末梢線は、逆行列を形成し、その後の計算を実施するための記憶デバイスの選択解除を示しうる。逆行列は、ガウス消去法または別の方法によって計算されてもよい。

30

【数 2 0】

$$\tilde{D},$$

の値が得られると、回復されるかまたは推定される任意のチェックサムベクトル

【数 2 1】

$$\tilde{C}$$

が、

【数 2 2】

$$\tilde{C} = F\tilde{D}$$

40

を使用して、データベクトル

【数 2 3】

$$\tilde{D},$$

に基づいて計算されてもよい。

【0 0 1 6】

並列 R S - R A I D データ記憶アーキテクチャは、各クラスタ内のデータおよびチェックサムを、全てのクラスタに転送されるかまたは配信される中間または部分和に集計しうる。中間データシンボル、中間チェックサムシンボル、データ記憶デバイスのクラスタへ

50

の割当てに関するクラスタ構成情報、およびデータ記憶デバイスの作動に関する状態などの使用は、並列RS-R A I D分散データ記憶アーキテクチャのスケラビリティおよびスループットを向上させながら、誤り訂正計算についての計算負荷およびレーテンシを低減しうる。

【0017】

本開示は、同じ数字が同じ要素を表す添付図面を参照することになる。

【図面の簡単な説明】

【0018】

【図1】並列RS-R A I D (Reed-Solomon redundant array of inexpensive disk)の例を示す図である。

10

【図2】構成行列の例を示す図である。

【図3】R A I Dコントローラの例を示す図である。

【図4A】チェックサムプログラムの例のフローチャートである。

【図4B】チェックサム更新プログラムの例のフローチャートである。

【図5】データプログラムの例のフローチャートである。

【発明を実施するための形態】

【0019】

図1は、データ記憶のための並列RS-R A I Dアーキテクチャ100の例である。並列RS-R A I Dアーキテクチャ100は、通信ファブリック1200、R A I Dコントローラ1111~1113、および記憶デバイス1001~1012を含みうる。記憶デバイス1001~1004、1005~1008、および1009~1012は、それぞれ、R A I Dコントローラ1111~1113に結合されうる。換言すれば、記憶デバイス1001~1012の部分集合またはクラスタは、それぞれの対応するR A I Dコントローラ1111~1113に結合しうる。各R A I Dコントローラ1111~1113に結合しうる記憶デバイスの数は、等しくても、等しくなくてもよく、また、耐故障性を改善する、スループットを改善するなどのために、記憶デバイスのR A I Dコントローラに対する構成またはマッピングは動的に変わってもよい。たとえば、記憶デバイス1001~1012のR A I Dコントローラ1111~1113に対する割当ては、構成行列または同様のデータ構造によって確定されてもよい。

20

【0020】

図2は、変数「t」を含みうる構成行列200の例を示し、変数「t」は、R A I Dコントローラの数のインデックスまたはカウンタである。たとえば、構成行列200の行206は、R A I Dコントローラインデックス番号「t」と、R A I Dコントローラ1111~1113などのそれぞれのR A I Dコントローラとの間のマッピング関数 $Q(t)$ を示す。行202はR A I D記憶デバイス開始インデックス $Q_S(t)$ を示し、行204はR A I D記憶デバイス終了インデックス $Q_E(t)$ を示す。たとえば、 $Q_S(2) = 1005$ であり、 $Q_E(2) = 1008$ である。デバイス番号のオフセットが、関数 $J(\cdot)$ によって供給されてもよいため、 $J(Q_S(2)) = 5$ であり、これは、たとえば、5番目の記憶デバイスが第2の記憶デバイスクラスタで開始することを示しうることに留意されてもよい。構成行列200は、記憶デバイスに対応するR A I Dコントローラにマッピングしうる。換言すれば、構成行列200は、記憶デバイスのどの部分集合またはクラスタが、所与のR A I Dコントローラに割当てられるかを確定するかまたは制御しうる。計算のために、構成行列200は、以下で述べるように、コードワードを符合化するかまたは復号化しうる重み付き部分和の始まりと終わりを確定しうる、チェックサムおよびデータを更新するかまたは維持しうるなどを行いうる。

30

40

【0021】

通信ファブリック1200は、R A I Dコントローラ1111~1113間で、また、並列RS-R A I Dアーキテクチャ100と外部デバイスとの間で入力および出力(I/O)デジタル信号を結合しうる。たとえば、通信ファブリック1200は、R A I Dコントローラ1111~1113間で、データシンボル、チェックサムシンボル、中間データ

50

およびチェックサムシンボルなどのようなデジタル信号を結合しうる。通信ファブリック 1200 は、並列バス構造、直列データリンク、光バックプレーンなどを使用してもよい。通信ファブリック 1200 は、外部通信に 1 つのタイプのバス、リンク、またはバックプレーン構造を、RAID コントローラ 1111 ~ 1113 間の通信に別のタイプを使用してもよい。

【0022】

RAID コントローラ 1111 ~ 1113 は、構成行列 200 などの構成行列または他のデータ構造によって与えられた、割当てられた記憶デバイスの部分集合またはクラスタ内の各記憶デバイスについてデータチェックサムシンボルを計算しうる。RAID コントローラ 1111 ~ 1113 は、誤り訂正コード計算の部分を集計するかまたは蓄積し、集計されたデータおよびパリティ計算結果を、通信ファブリック 1200 を通じて、並列 RS - RAID アーキテクチャ 100 内の他の RAID コントローラに報告しうる。データおよびチェックサムシンボルについての部分計算の詳細が、特定の RAID コントローラを参照して述べられてもよいが、対応する計算は、RAID コントローラ 1111 などの、並列 RS - RAID アーキテクチャ 100 内の任意の RAID コントローラによって実施されてもよい。

【0023】

図 3 は、通信ファブリックインタフェース 1111 a、RAID コントロールユニット 1111 b、中間和デバイス 1111 c、記憶デバイスインタフェース 1111 g、および記憶デバイス故障センスユニット 1111 h を含む RAID コントローラ 1111 の例を示す。通信ファブリックインタフェース 1111 a は、通信ファブリック 1200 などの通信ファブリックへのまた通信ファブリックからの信号を、中間和デバイス 1111 c および RAID コントロールユニット 1111 b に結合しうる。RAID コントロールユニット 1111 b は、中間和デバイス 1111 c、記憶デバイスインタフェース 1111 g、および記憶デバイス故障センスユニット 1111 h に結合しうる。記憶デバイスインタフェース 1111 g は、RAID コントロールユニット 1111 b、中間和デバイス 1111 c、および記憶デバイス故障センスユニット 1111 h に結合しうる。RAID コントローラ 1111 は、先に述べたように、通信ファブリック 1200 へまた通信ファブリック 1200 から結合し、記憶デバイスインタフェース 1111 g を介して記憶デバイス 1001 ~ 1004 などの記憶デバイスへまた記憶デバイスから結合しうる。

【0024】

中間和デバイス 1111 c は、中間和計算器 1111 d、再計算器 1111 e、および計算コントロール 1111 f を含む。中間和計算器 1111 d は、通信ファブリックインタフェース 1111 a、記憶デバイスインタフェース 1111 g、再計算器 1111 e、および計算コントロール 1111 f に結合しうる。再計算器 1111 e は、通信ファブリックインタフェース 1111 a、中間和計算器 1111 d、計算コントロール 1111 f、および記憶デバイスインタフェース 1111 g に結合しうる。計算コントロール 1111 f は、中間和計算器 1111 d、再計算器 1111 e、および記憶デバイスインタフェース 1111 g に結合しうる。

【0025】

通信ファブリックインタフェース 1111 a は、並列 RS - RAID アーキテクチャ 100 と外部デバイスとの間で情報シンボルを転送し、通信ファブリック 1200 と RAID コントローラ 1111 の要素との間で、情報シンボル、情報シンボルの所定部分、データシンボル、中間チェックサムシンボルなどのチェックサムシンボル、コントロール信号、クロック信号などを結合しうる。通信ファブリックインタフェース 1111 a は、情報シンボルをビットからバイト、ワード、または他のシンボルにリフォーマットしうる、信号を多重化し逆多重化しうる、データ転送を同期化しうる、ラインドライバおよび受信機によって信号をバッファリングしうるなどを行いうる。換言すれば、通信ファブリックインタフェース 1111 a は、デジタルバスなどの通信ファブリックを通じて送信するためにデジタル信号を調節しうる、データ転送をバッファリングしうるなどを行いうる。

【0026】

RAIDコントロールユニット1111bは、通信ファブリックインタフェース1111aおよび記憶デバイスから信号を受信しうる、情報シンボルの部分集合からデータシンボルを選択しうる、記憶デバイスにわたってデータおよびチェックサムシンボルをストライピングしうる、前方誤り訂正コード(forward-error correction code)(FEC code)によって中間和デバイス1111cの作動を制御しうるなどを行う。たとえば、情報シンボルの部分集合は、データシンボルによって表され、かつ、RAIDコントローラ1111によって制御される作動可能なデータ記憶デバイスに記憶される情報シンボルでありうる。中間和デバイス1111cは、記憶デバイス故障センスユニット1111hから状態情報を得ることができるRAIDコントロールユニット1111bから作動可能な記憶デバイスの数に関する状態情報を受信してもよい。

10

【0027】

記憶デバイス故障センスユニット1111hは、RAIDコントローラ1111に結合する任意の記憶デバイスの作動に関する状態を確定し、作動可能な記憶デバイスのリストを確定しうる。換言すれば、記憶デバイス故障センスユニット1111hは、所与の記憶デバイスが、データおよびチェックサムの確実な記憶に適さなくなったかどうかを判定しうる。記憶デバイス故障センスユニット1111hは、信頼性のある作動について記憶デバイスを試験しうる、所与の記憶デバイスがオンラインであるかどうかを判定しうる、所与の記憶デバイスからの応答が、所定のタイムアウト間隔内に受信されない場合、ユニットオフラインを宣言しうる、信号品質メトリックが、記憶デバイスから読出されたデータについて閾品質より小さいかどうかを判定しうる、作動可能な記憶デバイスを挙げうるなどを行う。記憶デバイス故障センスユニット1111hは、こうした試験の結果を記録し、RAIDコントロールユニット1111bなどのRAIDコントローラ1111の要素のために、作動可能な記憶デバイスのリストを配信しうる。

20

【0028】

中間和計算器1111dは、中間的でローカルで部分的な和を計算することができ、この和内に、チェックサムおよびデータについての誤り訂正コード計算が、式8および式13に関してそれぞれ述べたように分解されうる。中間的なまたは部分的な和は、RAIDコントローラ1111に報告するクラスタ内の作動可能な記憶デバイスから読出されるシンボルの重み付き和であってよい。たとえば、記憶デバイスのクラスタおよびこうした部分和の合計の対応する制限は、構成行列200などの構成行列または他のデータ構造から確定されてもよい。中間和計算器1111dは、RAIDコントローラ1112またはRAIDコントローラ1113などの他のRAIDコントローラから対応する部分和を受信した後、データおよびチェックサムシンボルを計算しうる。

30

【0029】

再計算器1111eは、RAIDコントローラ1111に直接結合する記憶デバイスからのデータに基づいて中間的なローカルのチェックサムを、また、通信ファブリックインタフェース1111aを通して転送される他のRAIDコントローラからのローカルでない中間チェックサムを再計算しうる。換言すれば、データまたはチェックサムシンボルの変化が、RAIDコントローラ1111に直接結合するローカルの記憶デバイスにおいて起こるか、または、通信ファブリック1200を通じてRAIDコントローラ1111に送信される中間チェックサムによって起こると、再計算器1111eは、中間和計算器1111dからの結果を相応して修正しうる。

40

【0030】

計算コントロール1111fは、中間チェックサム計算結果または再計算されたチェックサムが、FECのために使用されるべきかどうかを判定するために、中間和計算器1111dと再計算器1111eの両方を制御しうる。RAIDコントロール1111bは、中間和計算器1111dの結果と再計算器1111eの結果のいずれが計算されるかを判定するために、直接にまたは通信ファブリックインタフェース1111aを通して計算コ

50

ントロール 1 1 1 1 f に合図し(signal)うる。RAIDコントロール 1 1 1 1 b は、記憶デバイス故障センスユニット 1 1 1 1 h から、データ記憶デバイスに関する作動に関する状態などの状態情報を得ることができる。

【 0 0 3 1 】

並列 RAID コントローラ 1 1 1 1 ~ 1 1 1 3 は、

【 数 2 4 】

$$c_i = \sum_{t=1}^r \sum_{j=J(QS(t))}^{J(QE(t))} d_j f_{i,j} = \sum_{t=1}^r c_{i,t} \quad \text{式 7}$$

によって、チェックサムを計算し記憶しうる。式中、インデックス t は、1 から RAID コントローラの数 r までの範囲にあることができ、 $c_{i,t}$ は、 t 番目のインデックスについての i 番目の中間チェックサムである。たとえば、 r は、並列 RAID アーキテクチャ 1 0 0 の場合、3 に等しい。構成行列 2 0 0 に関して述べたように、 $QS(t)$ および $QE(t)$ は、開始および終了記憶デバイスを RAID コントローラにマッピングし、それぞれの中間チェックサム $c_{i,j}$ を生成する部分和の合計の制限を確定しうる。関数 $J(\cdot)$ は、たとえば、 $J(1 0 0 2) = 2$ であるように、オフセットを減算しうる。

10

【 0 0 3 2 】

RAID コントローラ 1 1 1 1 などの t 番目の RS - RAID コントローラは、

【 数 2 5 】

$$c_{i,t} = \sum_{j=J(QS(t))}^{J(QE(t))} d_j f_{i,j} \quad \text{式 8}$$

20

によって、中間チェックサム $c_{i,t}$ を計算しうる。

【 0 0 3 3 】

中間チェックサム $c_{i,t}$ の使用は、通信ファブリック 1 2 0 0 上のデータトラフィックを低減することができる。並列 RS - RAID アーキテクチャ 1 0 0 のスループットを増加させることができる。たとえば、8 + 4 RS - RAID アーキテクチャでは、単一の主要な RAID コントローラが、記憶デバイスの全てを制御し、チェックサムを計算する場合、8 つのデータシンボルが、通信ファブリックを通じて転送されうる。対照的に、8 + 4 並列 RS - RAID アーキテクチャからの中間チェックサム計算器結果を使用して、2 つの中間チェックサムシンボルだけが、通信ファブリックを通じて転送される必要がある可能性がある。

30

【 0 0 3 4 】

中間チェックサムおよび全チェックサムを計算することに加えて、並列 RS - RAID アーキテクチャ 1 0 0 は、データシンボルが変化すると、チェックサムシンボルを修正するかまたは維持しうる。たとえば、データシンボルが、 d_j から d'_j に変化すると、チェックサムは、

$$c'_i = c_i + f_{i,j} (d'_j - d_j) \quad \text{式 9}$$

によって、再計算されうる。式 9 の計算を実施するとき、RAID コントローラ 1 1 1 1 は、データ差 ($d'_j - d_j$) を計算し、ファンデルモンデ要素 $f_{i,j}$ 、すなわち、

40

【 数 2 6 】

$$c'_{i,t} = \sum_{j=J(QS(t))}^{J(QE(t))} f_{i,j} (d'_j - d_j) \quad \text{式 10}$$

によって、データ差に重み付けしうる。

【 0 0 3 5 】

個々の並列 RAID コントローラ 1 1 1 1 ~ 1 1 1 3 は、一時的な成分 $c'_{i,t}$ を RAID コントローラ 1 1 1 1 ~ 1 1 1 3 の他のコントローラに送出しうる。

【 0 0 3 6 】

RS - RAID コントローラ 1 1 1 1 ~ 1 1 1 3 は、

【数 2 7】

$$c'_i = c_i + \sum_{i=1}^r c'_{i,i} \quad \text{式 1 1}$$

によって、それぞれの割当てられた記憶デバイスを更新しうる。

【0 0 3 7】

記憶デバイスが故障する、たとえば、記憶デバイス故障センスユニット 1 1 1 1 h がハードディスククラッシュを検出すると、拡大行列の逆行列 $\text{Inv}(A)$ が、並列 RAID コントローラ 1 1 1 1 ~ 1 1 1 3 によって修正されて、残りのまたは作動可能なデータ記憶デバイスに対応する逆行列 $\text{Inv}(A')$ が形成される。行列 $\text{Inv}(A')$ は、記憶デバイスのさらなる故障が起こらない限り、静的データ構造である可能性がある。別の記憶デバイスが故障すると、 $\text{Inv}(A')$ が一回計算され、その後、RAID コントローラ 1 1 1 1 ~ 1 1 1 3 などの全ての作動可能な RAID コントローラにブロードキャストされうる。より多くの記憶デバイスが後で故障する場合、新しい逆行列 $\text{Inv}(A'')$ が、再計算され、全ての RAID コントローラにブロードキャストされてもよい。

10

【0 0 3 8】

並列 RS - RAID アーキテクチャ 1 0 0 は、記憶デバイスが故障しても、各 RAID コントローラにおいてローカルで計算される中間和または部分和を使用して、データシンボルを回復しうる。回復されるデータ

【数 2 8】

$$\tilde{D} = [\tilde{d}_1, \tilde{d}_2, \dots, \tilde{d}_n]^T$$

20

は、

【数 2 9】

$$\begin{bmatrix} \tilde{d}_1 \\ \tilde{d}_2 \\ \vdots \\ \tilde{d}_n \end{bmatrix} = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,n} \\ a_{2,1} & a_{2,2} & \dots & a_{2,n} \\ \vdots & \vdots & & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,n} \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad \text{式 1 2}$$

から回復される可能性がある。式中、 $\text{Inv}(A')$ の要素は、 $a_{i,j}$ ($1 \leq i \leq n$ および $1 \leq j \leq n$) で示されてもよい。対応するデータおよびチェックサムシンボル $E' = [e_1, e_2, \dots, e_n]^T$ の要素は、作動可能でかつ選択されたデータ記憶デバイスから読出されうる。並列 RS - RAID アーキテクチャ 1 0 0 は、E の行ならびに追加された単位行列とファンデルモンデ行列の対応する部分集合を選択するかまたは選抜して、それぞれ、 E' および $\text{Inv}(A')$ が形成されうる。換言すれば、並列 RS - RAID アーキテクチャ 1 0 0 は、データ回復計算を、

30

【数 3 0】

$$\tilde{d}_{i,t} = \sum_{j=J(QS(t))}^{J(QE(t))} a_{i,j} \cdot e_j \quad \text{式 1 3}$$

40

によって、部分和または中間データシンボルの集合に分解しうる。式中、 e_j は、t 番目の RS - RAID コントローラの制御下にある全てのデータまたはチェックサムシンボルの集合であると理解される。

【0 0 3 9】

他の並列 RS - RAID コントローラから、中間データシンボルなどのメッセージを受信すると、個々の RAID コントローラは、最初に中間データシンボル

【数 3 1】

$$\tilde{d}_{i,t}$$

を計算し、次に、

50

【数 3 2】

$$\tilde{d}_i = \sum_{t=1}^r \tilde{d}_{i,t}. \quad \text{式 1 4}$$

によって、回復されるデータ

【数 3 3】

$$\tilde{d}_i$$

を計算しうる。

【0040】

図 4 A は、データ記憶用の並列 R S - R A I D アーキテクチャのためのチェックサムプログラムフローチャート 4 0 0 A の例を示す。プログラムフローチャート 4 0 0 A は、プログラムステップ S 4 1 0 で開始し、並列 R S - R A I D アーキテクチャの構成行列が読出されうるプログラムステップ S 4 2 0 に進みうる。たとえば、構成行列は、図 2 に関して述べた構成行列などの、所与の R A I D コントローラに関連する記憶デバイスについての開始デバイス番号と終了デバイス番号を指定しうる。各 R A I D コントローラは、構成行列のローカルなコピーを記憶しうる、構成行列を他の R A I D コントローラと調和させうる、構成行列を、高レベルの R A I D デバイスまたはネットワークコントローラから受信しうるなどを行いうることが理解されてもよい。

10

【0041】

プログラムステップ S 4 2 0 から、プログラムフローは、プログラムステップ S 4 2 5 に進むことができ、プログラムステップ S 4 2 5 にて、プログラムは、外部デバイスから、記憶される情報保持データシンボルを読出すことができる。たとえば、プログラムは、通信ファブリックを通じて受信される 2 K ビットのデータブロックの集合をフラッシュドライブから受信しうる。

20

【0042】

プログラムステップ S 4 2 5 から、プログラムは、中間チェックサムが計算されうるプログラムステップ S 4 3 0 に進むことができる。たとえば、中間チェックサムまたはデータおよびパリティ計算は、式 8 およびガロア体の特性を使用した、データワードの線形結合から計算されうる。プログラムステップ S 4 3 0 は、中間チェックサムを計算し、1) 所与の R A I D コントローラにコードワードシンボルの所定部分を供給する個々の作動可能な記憶ユニットからの記憶されたコードワードシンボルを使用して中間チェックサムを更新するかまたは維持し、2) 通信ファブリックを通じて所与の R A I D コントローラと通信する他の R A I D コントローラからの中間チェックサムを集計しうる。換言すれば、ローカルな部分集合データ記憶ユニットからの低レテンシのデータおよびパリティビット、バイト、またはワードは、他のデータ記憶ユニットからの、 $c_{i,t}$ の形態の、高レテンシの、蓄積されるかまたは部分的に合計されたデータおよびパリティと結合されうる。その後、プログラムは、プログラムステップ S 4 6 0 に進む。

30

【0043】

ステップ S 4 6 0 にて、プログラムは、中間チェックサムを異なる R S - R A I D コントローラに配信しうる。たとえば、 $Q(t) = t$ である場合、プログラムステップ S 4 6 0 は、第 1 の R A I D コントローラからの第 1 の中間チェックサム $c_{1,1}$ を第 2 および第 3 の R A I D コントローラに配信しうる。

40

【0044】

プログラムステップ S 4 6 0 から、プログラムフローは、プログラムステップ S 4 7 0 に進むことができ、プログラムステップ S 4 7 0 にて、プログラムは、他の R A I D コントローラから中間チェックサムを受信しうる。プログラムステップ S 4 7 0 から、プログラムはプログラムステップ S 4 8 0 に進みうる。中間チェックサムの集合によって、各 R A I D コントローラが、式 8 によって完全なチェックサム c_i を計算し、その後の誤り訂正および検出計算について c_i を記憶することが可能になる。たとえば、プログラムは、第 2 および第 3 の中間チェックサム $c_{i,2}$ および $c_{i,3}$ を受信することができ、第 2

50

および第3の中間チェックサム $c_{i,2}$ および $c_{i,3}$ は、ローカルに計算された第1のチェックサム $c_{i,1}$ と共に、 $c_{i,1}$ を計算するためのチェックサムの十分な集合を形成しうる。

【0045】

プログラムステップS480から、プログラムフローは、プログラムステップS490に進むことができ、プログラムステップS490にて、プログラムは、プログラムを実行するRAIDコントローラに割当てられるデータおよび完全なチェックサムシンボルを記憶しうる。たとえば、プログラムは、ディスクのアレイにわたってデータおよびチェックサムシンボルをストライピングしうる。プログラムステップS490から、プログラムフローは、プログラム実行が停止しうるプログラムステップS495に進みうる。

10

【0046】

図4Bは、データ記憶用の並列RS-Raidアーキテクチャのためのチェックサム更新プログラムフローチャート400Bの例を示す。プログラムフローチャート400Bは、ステップS440で開始し、ステップS442に進む。

【0047】

ステップS442にて、並列RS-Raidアーキテクチャは、データの変化を受信する可能性がある。たとえば、記憶デバイスは、古いデータシンボルを置換するために新しいデータシンボルを受信してもよい。その後、プログラムフローはステップS444に進みうる。

【0048】

ステップS444にて、記憶デバイスに結合するRAIDコントローラは、式10によって、一時的な成分を計算しうる。RAIDコントローラは、新しいデータシンボルと古いデータシンボルとのデータ差を得、ファンデルモンデ行列要素によってデータ差に重み付けしてもよい。その後、プログラムフローはステップS446に進みうる。

20

【0049】

ステップS446にて、一時的な成分は、他のRAIDコントローラに伝達されうる。ある実施形態では、通信ファブリックは、種々のRAIDコントローラを結合してもよい。通信ファブリックは、データ変化に対応する一時的な成分を、チェックサムを記憶する記憶デバイスを制御するRAIDコントローラに伝達しうる。その後、プログラムフローはステップS448に進みうる。

30

【0050】

ステップS448にて、チェックサムを記憶する記憶デバイスを制御するRAIDコントローラは、たとえば、式11によって、受信した一時的な成分に基づいてチェックサムを更新してもよい。その後、プログラムフローはステップS450に進み、停止しうる。

【0051】

図5は、データ記憶用の並列RS-Raidアーキテクチャのためのデータプログラムフローチャート500の例を示す。プログラムフローチャート500は、ステップS510で開始し、ステップS520に進むことができ、ステップS520にて、並列RS-Raidアーキテクチャの構成行列が、図4Aに関して説明したように読出されうる。プログラムステップS520から、プログラムフローは、プログラムステップS525に進むことができ、プログラムステップS525にて、データおよびチェックサムシンボルが、記憶デバイスから読出されうる。たとえば、8つのデータおよび4つのチェックサムシンボルは、12の記憶デバイスから読出されうる。この例では、少なくとも8つのデータまたはチェックサムシンボルが、作動可能な記憶デバイスから読出されうる。

40

【0052】

プログラムステップS425から、プログラムフローは、プログラムステップS530に進むことができ、プログラムステップS530にて、プログラムは、中間データシンボルを計算しうる。たとえば、プログラムは、式13によって、中間データシンボルを計算しうる。式13で使用される重み係数 $a_{i,j}$ は、予め計算され、RAIDコントローラに配信されるか、または、プログラムステップS520において構成行列を読出した後な

50

どに、必要に応じて再計算されてもよいことが理解されてもよい。プログラムステップ S 5 3 0 から、プログラムフローは、プログラムステップ S 5 4 0 に進むことができ、プログラムステップ S 5 4 0 にて、プログラムは、中間データシンボルを並列 R A I D コントローラに配信しうる。

【 0 0 5 3 】

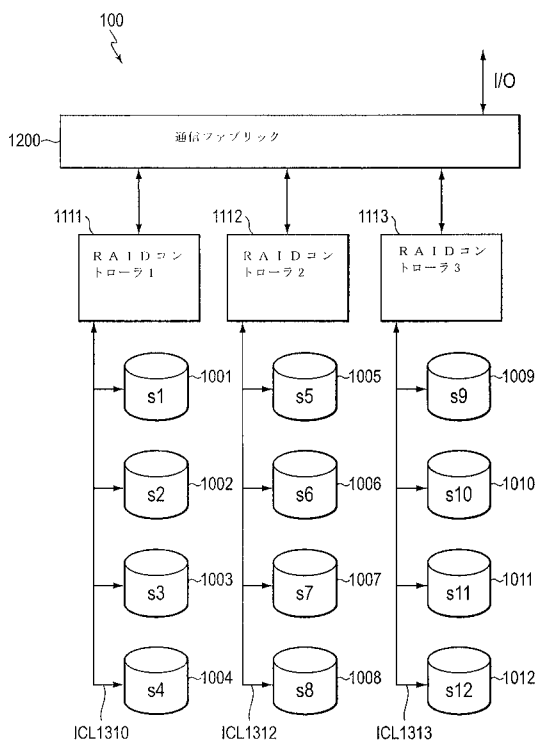
プログラムステップ S 5 4 0 から、プログラムフローは、プログラムステップ S 5 5 0 に進むことができ、プログラムステップ S 5 5 0 にて、プログラムは、並列 R A I D コントローラから中間データシンボルを受信しうる。プログラムステップ S 5 5 0 から、プログラムフローは、プログラムステップ S 5 6 0 に進むことができ、プログラムステップ S 5 6 0 にて、プログラムは、ローカルな R A I D コントローラと並列 R A I D コントローラの両方から来た中間データシンボルから、回復されるデータシンボルを計算しうる。換言すれば、プログラムは、式 1 4 によって、中間データシンボルを合計しうる。プログラムステップ S 5 6 0 から、プログラムフローは、プログラム実行が停止しうるプログラムステップ S 5 7 0 に進むことができる。

10

【 0 0 5 4 】

本発明は、本発明の特定の例示的な実施形態に関連して述べられたが、多くの代替、修正、および変形が当業者に明らかになることが明らかである。したがって、本明細書で述べる本発明の実施形態は、制限的でなく、例証的であることを意図される。本発明の精神および範囲から逸脱することなく行われてもよい変更が存在する。

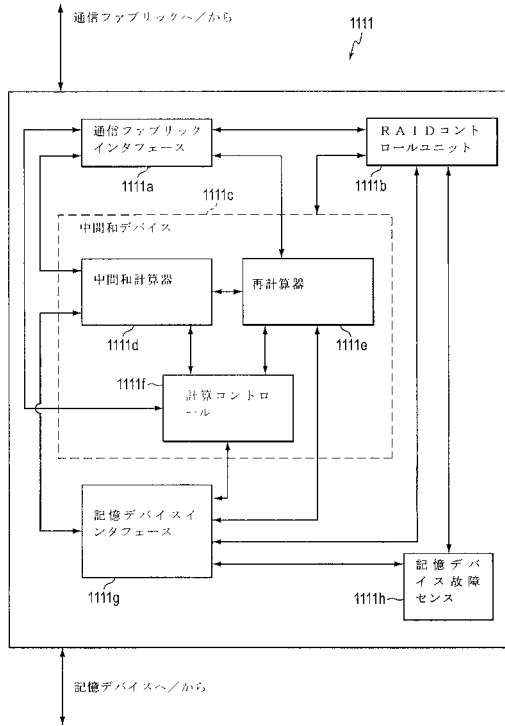
【 図 1 】



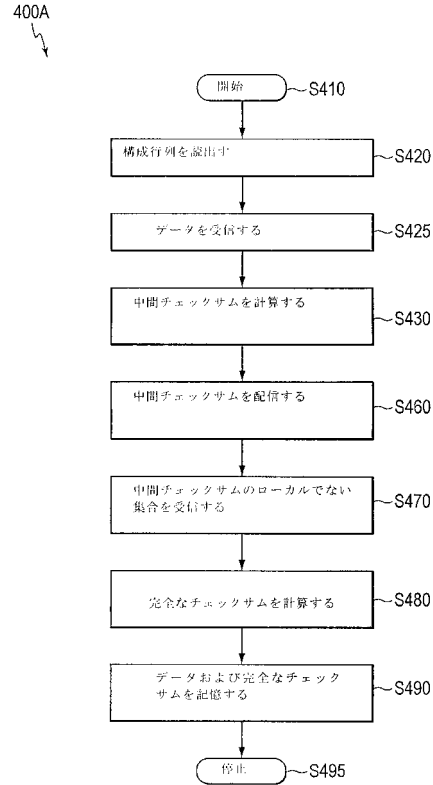
【 図 2 】

200				
	t	1	2	3
202	QS (t)	1001	1005	1009
204	QE (t)	1004	1008	1012
206	Q (t)	1111	1112	1113

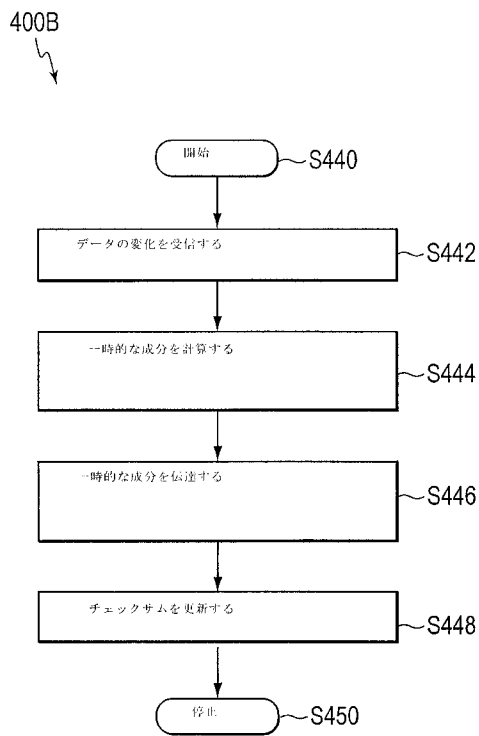
【図3】



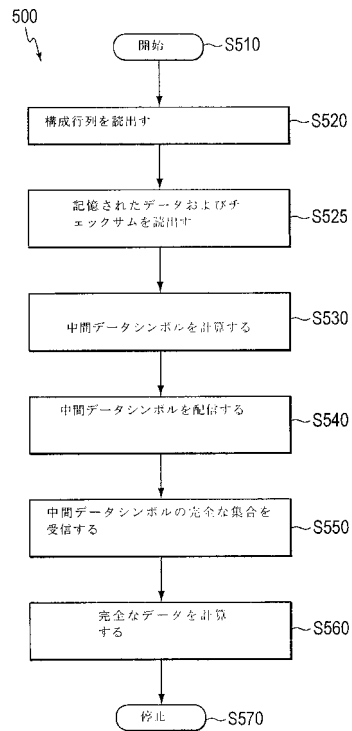
【図4A】



【図4B】



【図5】



フロントページの続き

(56)参考文献 特開2007-257630(JP,A)
特開2003-131818(JP,A)
米国特許出願公開第2007/0245173(US,A1)

(58)調査した分野(Int.Cl., DB名)
G06F 3/06