

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7187387号

(P7187387)

(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 11/10 (2006.01)

G 0 6 F 11/10 6 0 4

G 0 6 F 3/06 (2006.01)

G 0 6 F 3/06 3 0 4 B

G 0 6 F 13/10 (2006.01)

G 0 6 F 3/06 3 0 5 C

G 0 6 F 3/06 3 0 2 J

G 0 6 F 13/10 3 4 0 A

請求項の数 19 (全24頁)

(21)出願番号 特願2019-103825(P2019-103825)

(22)出願日 令和1年6月3日(2019.6.3)

(65)公開番号 特開2019-212310(P2019-212310
A)

(43)公開日 令和1年12月12日(2019.12.12)

審査請求日 令和4年5月20日(2022.5.20)

(31)優先権主張番号 62/682,763

(32)優先日 平成30年6月8日(2018.6.8)

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 16/103,907

(32)優先日 平成30年8月14日(2018.8.14)

(33)優先権主張国・地域又は機関

米国(US)

早期審査対象出願

(73)特許権者 390019839

三星電子株式会社

Samsung Electronics
Co., Ltd.大韓民国京畿道水原市靈通区三星路12
9129, Samsung-ro, Yeon-
gtong-gu, Suwon-si
, Gyeonggi-do, Repub-
lic of Korea

(74)代理人 100107766

弁理士 伊東 忠重

(74)代理人 100070150

弁理士 伊東 忠彦

(74)代理人 100091214

最終頁に続く

(54)【発明の名称】 低帯域データリペアを補助するシステム、装置及び方法

(57)【特許請求の範囲】

【請求項1】

装置であって、

データエラー訂正のための少なくとも一つのタイプのデータ再生成コードを演算するよ
うに構成された再生成コード認識(RCA)ストレージ装置、

を含み、

前記RCAストレージ装置は、

データブロックを含む複数のチャンクにてデータを格納するように構成されたメモリ
と、外部のホスト装置と関連付けられた要請に基づき、選択された数のデータブロックに
基づいてデータ再生成コードを演算するように構成されたプロセッサと、

インターフェースであり、

前記データ再生成コードを前記外部のホスト装置に伝送し、

前記データ再生成コードを演算するように前記プロセッサを構成するコマンドを、
前記外部のホスト装置から受信する、

ように構成されたインターフェースと、

を含む、

装置。

【請求項2】

前記RCAストレージ装置は、

10

20

互いに異なるデータ再生コードを生成するように構成された複数の命令セットを格納するように構成されたコードメモリ、

をさらに含み、

前記プロセッサは、前記外部のホスト装置によって、前記データ再生コードを演算するために前記複数の命令セットのうちの一つを選択するように構成される、

請求項 1 に記載の装置。

【請求項 3】

前記コードメモリは、前記複数の命令セットが前記外部のホスト装置によって前記コードメモリに書き込まれるように構成される、請求項 2 に記載の装置。

【請求項 4】

前記インターフェースは、選択された数のデータブロックに基づく前記データ再生コードの生成を可能にするコマンドを、前記外部のホスト装置から受信するように構成される、請求項 1 に記載の装置。

【請求項 5】

前記インターフェースは、

リペアデータが要求されることを示すとともに前記データ再生コードが演算されるべきであることを示すリペアコマンドを前記外部のホスト装置から受信し、

前記データ再生コードを前記外部のホスト装置に返す、

ように構成され、

前記データ再生コードのサイズは、前記選択された数のデータブロックのサイズよりも小さい、

請求項 1 に記載の装置。

【請求項 6】

前記プロセッサは、前記外部のホスト装置によって要請されたときに、データ再生技法を通して前記データ再生コードの異なるバージョンを演算するように構成され、

前記プロセッサによって演算されるバージョンは、前記外部のホスト装置によって決定される、

請求項 1 に記載の装置。

【請求項 7】

システムであって、

ホスト装置であり、

分散ストレージシステム間で複数のデータチャンクとしてデータを格納し、

あるデータチャンクがエラーと関連付けられたことを検出し、

前記エラーの前記検出に応答して、データ再生技法を通して、前記複数のデータチャンクに基づいて、前記エラーと関連付けられた前記データチャンクを再構成する、

ように構成されたホスト装置と、

前記分散ストレージシステムであり、

それぞれのデータチャンクを格納するように構成された複数のストレージ装置を含み、

前記複数のストレージ装置は、少なくとも一つのタイプのデータ再生コードを内部で演算するように構成された少なくとも一つの再生コード認識 (RCA) ストレージ装置を含む、

前記分散ストレージシステムと、

を含むシステム。

【請求項 8】

前記 RCA ストレージ装置は、

データブロックを含むチャンクにてデータを格納するように構成されたメモリと、

選択された数のデータブロックに基づいてデータ再生コードを演算するように構成されたプロセッサと、

前記データ再生コードを前記ホスト装置に伝送するように構成された外部インターフェースと、

10

20

30

40

50

を含む、請求項 7 に記載のシステム。

【請求項 9】

前記ホスト装置は、
データ再生コードを内部で演算できるストレージ装置を決定し、
該ストレージ装置からデータのチャンク又はその一部を要請し、
前記データのチャンク又はその一部に少なくとも部分的に基づいて、データ再生コードを前記ホスト装置により演算する、
ように構成される、請求項 7 に記載のシステム。

【請求項 10】

前記ホスト装置は、前記データ再生コードの前記演算を前記ストレージ装置にオフロードすることを、以下のファクタうちの一つ以上、すなわち、
前記ストレージ装置に利用可能なデータ再生技法、
前記分散ストレージシステムに関連付けられた利用可能な帯域幅の量、
前記データのチャンク又はその一部のサイズと比較した前記データ再生コードのサイズ、及び
前記ホスト装置内の利用可能な計算能力の大きさ、
のうちの一つ以上に少なくとも部分的に基づいて決定するように構成される、請求項 9 に記載のシステム。

【請求項 11】

前記ホスト装置は、前記ホスト装置によって演算されたデータ再生コードとそれぞれのストレージ装置によって演算された前記データ再生コードとに基づいて、前記エラーと関連付けられた前記チャンクを再構成するように構成される、請求項 9 に記載のシステム。

【請求項 12】

前記ホスト装置は、
データ再生コードを内部で演算できる第 1 のストレージ装置と第 1 プロトコルを介して通信し、
データ再生コードを内部で演算できない第 2 のストレージ装置と第 2 プロトコルを介して通信する、
ように構成される、請求項 9 に記載のシステム。

【請求項 13】

前記ホスト装置は、
それぞれのデータ再生コードを内部で演算できるストレージ装置を検出し、
データ再生技法と関連付けられた命令を、該ストレージ装置に格納することで、該ストレージ装置が、該データ再生技法を通して前記データ再生コードを演算するように構成されるようにする、
ように構成される、請求項 7 に記載のシステム。

【請求項 14】

前記ホスト装置は、
それぞれのデータ再生コードを内部で演算できるストレージ装置を、少なくとも部分的に、前記ホスト装置によって選択されたデータ再生技法を通して前記データ再生コードを演算できるストレージ装置を検出することによって検出する、
ように構成される、請求項 13 に記載のシステム。

【請求項 15】

システムであって、
ホスト装置であり、
ストレージシステムの間で複数のチャンクにてデータを格納し、
あるチャンクがエラーと関連付けられたことを検出し、
前記エラーの前記検出に応答して、データ再生技法を通して、前記複数のチャンクに少なくとも部分的に基づいて、前記エラーを訂正する、

10

20

30

40

50

ように構成されたホスト装置と、
前記ストレージシステムであり、

前記データのそれぞれのチャンクを格納するように構成された複数のストレージ装置、
を含む前記ストレージシステムと、
を含み、

前記複数のストレージ装置は、少なくとも一つのタイプのデータ再生コードを内部で
演算するように構成された少なくとも一つの再生コード認識（RCA）ストレージ装置
を含み、

前記RCAストレージ装置は、

データブロックを含む複数のチャンクにてデータを格納するように構成されたメモリ
と、

前記ホスト装置と関連付けられた要請に基づき、選択された数のデータブロックに基
づいてデータ再生コードを演算するように構成されたプロセッサと、
互いに異なるデータ再生コードを生成するように構成された複数の命令セットを格納す
るように構成されたコードメモリと、

前記データ再生コードを前記ホスト装置に伝送するように構成された外部インター
フェースと、

を含む、

システム。

【請求項 16】

前記ホスト装置は、命令セットを前記RCAストレージ装置の前記コードメモリに書き
込むように構成され、該命令セットは、前記データ再生技法を通しての前記RCAスト
レージ装置による前記演算を可能にするように構成される、請求項 15 に記載のシステム。

【請求項 17】

前記ホスト装置は、それぞれのRCAストレージ装置に一つ以上のデータ再生コード
の前記演算を少なくとも部分的に動的にオフロードすることにより、前記エラーを訂正す
るように構成される、請求項 15 に記載のシステム。

【請求項 18】

前記ホスト装置は、前記ストレージシステム内のRCAストレージ装置ではないストレ
ージ装置によって格納されたデータの一つ以上のチャンクについてのデータ再生コードを
、前記ホスト装置によって演算することによって、前記エラーを訂正するように構成され、
前記ホスト装置によって演算することは、前記ストレージ装置から前記データのチャン
クの少なくとも一部を伝送することを含み、

前記RCAストレージ装置によって演算され、前記ホスト装置に伝送される前記データ
再生コードのサイズは、前記ホスト装置に伝送される前記ストレージ装置からの前記デ
ータのチャンクの前記少なくとも一部のサイズよりも小さい、請求項 17 に記載のシステ
ム。

【請求項 19】

前記複数のストレージ装置は、非RCA（non-RCA）ストレージ装置を含み、
前記ホスト装置は、

第1プロトコルを介して前記RCAストレージ装置と通信し、

第2プロトコルを介して前記非RCAストレージ装置と通信する、

ように構成される、

請求項 15 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はデータストレージに関するもので、より詳しくは低帯域データリペアを補助す
るストレージ装置のためのシステム、装置及び/又は方法に関する。

【背景技術】

10

20

30

40

50

【 0 0 0 2 】

コーディング理論 (coding theory) において、イレイジャーコード (erasure code) は、オリジナルメッセージが n 個のシンボルの一部から回復されることができるよう k 個のシンボルのメッセージを n 個のシンボルのより長いメッセージ (コードワード (code word)) に変換するビットイレイジャー (ビットエラーの代わりに) の仮定下での前方誤り訂正 (FEC: forward error correction) コードである。 $r=k/n$ なる比は、符号レート (code rate) と呼ばれる。 $(k - /k)$ なる比 (ここで、 $k -$ は回復のために要求されるシンボルの個数を示す) は、受信効率 (reception efficiency) と呼ばれる。

【 0 0 0 3 】

再生成コードは、既存のエンコードされたフラグメント (fragments) から、損失されたエンコードされたフラグメントを再構成 (又はリペア (repairing) と呼ばれる) するイシューを解決する。より詳細には、再生成コードは、従来の MDS (maximum distance separable) コードのストレージ効率性を維持しつつ、リペアにおけるダウンロードのサイズを減らすことを目標とするコードのクラスである。斯かるイシューは、エンコードされた冗長性を維持するための通信が問題になる分散ストレージシステムにおいて発生する。

10

【 0 0 0 4 】

分散ストレージシステム (distributed storage system) は一般的に、情報が、しばしば複製方式で、2 つ以上のノード又は装置に格納されるコンピュータネットワークである。分散ストレージシステムは、ユーザーが多数のノードに情報を格納する分散データベース又はユーザーが多数のピアネットワークノードに情報を格納するコンピュータネットワークのいずれかを表すのに使用される。分散ストレージシステムは、一般的に、エラー検出及び訂正技法を使用する。一部の分散ストレージシステムは、ファイルの一部が損傷されたり、又は利用可能ではなかったりする場合、前方誤り訂正技法 (forward error correction techniques) を使用してオリジナルファイル、チャンク又はバイナリラージオブジェクト (blob) を回復する。他の分散ストレージシステムは、他のミラーからファイルをダウンロードしようと再び試みる。

20

【 先行技術文献 】

【 特許文献 】

【 0 0 0 5 】

【 文献 】 米国特許第 9 7 8 5 4 9 8 号明細書

30

米国特許出願公開第 2 0 1 7 0 3 0 8 4 3 7 号明細書

【 非特許文献 】

【 0 0 0 6 】

【 文献 】 CALIS, Gokhan, "Coding and Maintenance Strategies for Cloud Storage: Correlated Failures, Mobility and Architecture Awareness," text; Electronic Dissertation, The University of Arizona, 2017, found via Google Scholar (url: http://arizona.openrepository.com/arizona/bitstream/10150/625607/1/azu_etd_15691_sip1_m.pdf), 164 pages.

RASHMI, K.V., et al., "Regenerating Codes for Errors and Erasures in Distributed Storage," IEEE International Symposium on Information Theory (ISIT), 2012, 5 pages.

40

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 7 】

本発明の目的は、本発明の実施形態によると、向上された信頼性及び改善された性能を有する低帯域データリペアを補助するシステム、装置及び方法を提供することにある。

【 課題を解決するための手段 】

【 0 0 0 8 】

本発明の一実施形態によると、装置は、データエラーを訂正するための少なくとも一つのタイプのデータ再生成コードを演算するように構成された再生成コード認識 (RCA: re

50

generation-code-aware)ストレージ装置を含むことができる。前記RCAストレージ装置は、それぞれがデータブロックを含むチャンクにてデータを格納するように構成されたメモリを含むことができる。前記RCAストレージ装置は、外部のホスト装置によって要請された場合、選択された数のデータブロックに基づいてデータ再生コードを演算するように構成されたプロセッサを含むことができる。前記RCAストレージ装置は、前記データ再生コードを前記外部のホスト装置に伝送するように構成された外部インターフェースを含むことができる。

【0009】

本発明の一実施形態によると、システムは、ホスト装置及び分散ストレージシステムを含むことができる。前記ホスト装置は、前記分散ストレージシステム間で複数のチャンクとしてデータを格納し、少なくとも一つのチャンクがエラーと関連付けられた場合を検出し、前記エラーに応答し、データ再生技法を通して、少なくとも前記データの前記複数のチャンクに基づいて、前記エラーと関連付けられた前記少なくとも一つのチャンクを再構成することができる。前記分散ストレージシステムは、複数のストレージ装置を含むことができる。前記複数のストレージ装置のそれぞれは、少なくとも前記データのそれぞれのチャンクを格納するように構成されることができる。前記複数のストレージ装置は、少なくとも一つの再生コード認識(RCA:regeneration-code-aware)ストレージ装置を含み、前記少なくとも一つのRCAストレージ装置のそれぞれは、少なくとも一つのタイプのデータ再生コードを内部で演算するように構成される。

【0010】

本発明の一実施形態によると、システムは、ホスト装置及びストレージシステムを含むことができる。前記ホスト装置は、前記ストレージシステムの間で複数のチャンクにてデータを格納し、少なくとも一つのチャンクがエラーと関連付けられていることを検出し、前記エラーに応答し、データ再生技法を通して前記データの前記複数のチャンクに基づいて、前記エラーを訂正するように構成されることができる。前記ストレージシステムは、複数のストレージ装置を含むことができる。前記複数のストレージ装置のそれぞれは、前記データの少なくともそれぞれのチャンクを格納するように構成され、前記複数のストレージ装置は、少なくとも一つの再生コード認識(RCA:regeneration-code-aware)ストレージ装置を含むことができる。前記RCAストレージ装置のそれぞれは、少なくとも一つのタイプのデータ再生コードを内部で演算するように構成されることができる。前記RCAストレージ装置は、チャンクにてデータを格納するように構成されたメモリを含むことができる。前記チャンクのそれぞれは、データブロックを含むことができる。前記RCAストレージ装置は、前記ホスト装置によって要請された場合、選択された数のデータブロックに基づいてデータ再生コードを演算するように構成されたプロセッサを含むことができる。前記RCAストレージ装置は、複数のコマンドのセットを格納するように構成されたコードメモリを含むことができる。前記複数のコマンドのセットのそれぞれは、他のデータ再生コードを生成することができる。前記RCAストレージ装置は、前記データ再生コードを前記ホスト装置に伝送するように構成された外部インターフェースを含むことができる。

【0011】

一つ又はそれ以上の具現に対する詳細な説明は、添付された図面及び以下の詳細な説明から掲示される。他の特徴は詳細な説明、図面及び特許請求の範囲から明確になるであろう。

【0012】

実質的に図面の中での少なくとも一つと関連付けられて説明されたり、又は図示され、特許請求の範囲でより完全に示されるように、本発明は、データストレージのためのシステム及び/又は方法、より詳しくは、ストレージ装置による低帯域幅のデータリペアに関するものである。

【発明の効果】

【0013】

本発明の実施形態によると、向上された信頼性及び改善された性能を有する低帯域データリペアを補助するシステム、装置及び方法が提供される。

【図面の簡単な説明】

【0014】

【図1】図1は、本発明の実施形態によるシステムの例示的な実施形態のブロック図である。

【図2a】図2aは、本発明の実施形態によるシステムの例示的な実施形態のブロック図である。

【図2b】図2bは、本発明の実施形態によるシステムの例示的な実施形態のブロック図である。

【図2c】図2cは、本発明の実施形態によるシステムの例示的な実施形態のブロック図である。

【図2d】図2dは、本発明の実施形態によるシステムの例示的な実施形態のブロック図である。

【図3】図3は、本発明の実施形態による技法の例示的な実施形態のフローチャートである。

【図4】図4は、本発明の実施形態の原理に基づいて形成された装置を含むことができる情報処理システムの概念的なブロック図である。

【0015】

多様な図面において、類似の参照番号は類似の要素を示す。

【発明を実施するための形態】

【0016】

多様な例示的な実施形態が、一部の例示的な実施形態を示す添付図面を参照して、さらに詳細に説明される。しかし、本文に開示された内容は、他の多様な形態で具現されることができ、本文に説明された例示的な実施形態に制限されると解釈されてはならない。代わりに、このような例示的な実施形態は、詳細な説明が完全であり、当業者によって、本発明の技術的思想が完全に伝達されるように提供される。図面で、階層と領域のサイズ及び相対的なサイズは、明確性のために、誇張されている。

【0017】

構成又は階層が、他の構成又は階層に「連結された(on、connected to、or coupled to)」ものと称される場合、これは直接的に(directly)他の構成又は階層と連結されることができたり、又は一つ以上の介在(intervening)構成又は階層が存在できることが理解されるだろう。対照的に、構成が他の構成又は階層と「直接的に連結された(directly on、directly connected to、or directly coupled to)」ものと称される場合には、介在構成又は階層が存在しない。類似の参照番号は、全体的に類似の構成を指し示す。本文で使用されるように、「及び/又は(and/or)」は、関連付けられて羅列されたリストのいずれか一つ、又はそれ以上の組み合わせの一部及び全部を含む。

【0018】

「第1(first)」、「第2(second)」、「第3(third)」などのような用語は、多様な要素、構成、領域、階層、及び/又はセクションを説明するために本文で使用されるが、このような要素、構成、領域、階層、及び/又はセクションは、斯かる用語に限定されないことが、よく理解されるだろう。斯かる用語は、一つの要素、構成、領域、階層、又はセクションを他の一つの要素、構成、領域、階層、又はセクションから区別するためにだけ使用される。つまり、以下で記載されている第1の要素、構成、領域、階層、又はセクションは、本発明の思想及び範囲から逸脱せずに、第2の要素、構成、領域、階層、又はセクションと称されることができる。

【0019】

「～の下(beneath、below、lower、under)」、「～の上(above、upper)」などのような空間的に相対的な用語(spatially relative terms)は、図面に図示した他の一つの要素(複数の要素)又は特徴(複数の特徴)と1つの要素又は特徴との関連性を容易に説明す

10

20

30

40

50

るために本文で使用されることができる。空間的に相対的な用語は、図面に図示された向きに加えて動作又は使用において装置の他の向きを含むものとして意図されていることがよく理解されるだろう。例えば、図面で、装置が覆される場合、他の要素又は特徴の「下(below or beneath or under)」に説明された要素は、他の要素又は特徴の「上(above)」に存在することになる。つまり、「下(below、under)」の例示的な用語は、上と下の両方の方向を含むことができる。装置は別の方向(例えば、90度回転するか、又は他の方向)を向くことができ、本文で使用される空間的に相対的な説明は、それに応じて解釈されるべきである。

【0020】

同様に、「ハイ(high)」、「ロー(low)」、「プルアップ(pull up)」、「プルダウン(pull down)」、「1」及び「0」などのような電氣的な用語は図面で示されたように、他の電圧レベル、他の構成又は特徴との相対的な電圧レベル又は電流を示す説明の便宜のために詳細な説明で使用されることができる。電氣的に相対的な用語は、図面に示された電圧又は電流に加えて使用又は動作で、装置の他の基準電圧を含むものとして意図されている。例えば、図面で、装置又は信号が反転されたり、他の基準電圧、電流、又は電荷を使用したりする場合、「ハイ(high)」又は「プルアップ(pull up)」で説明されている構成は、新しい基準電圧又は電流と比較して、「ロー(low)」又は「プルダウン(pull down)」であり得る。つまり、「ハイ(high)」の例示的な用語は、比較的に低い又は高い電圧及び電流の両方を含むことができる。装置は、他の電氣的なフレームの基準に基づくことができ、詳細な説明で使用される電氣的に相対的な説明により解釈されることができる。

【0021】

詳細な説明で使用される用語は、ただ特定の実施形態についての説明の目的のためのものであり、本発明の限定を意図するものではない。詳細な説明で使用されるように、単数形(singular forms)は、明確に別な方法で定義されていない限り、複数形(plural forms)を含むように意図される。「含む(comprise)」という用語は、詳細な説明に使用される場合、列挙された特徴、ステップ、動作、要素、及び/又は構成の存在を特定するものの、一つ又はそれ以上の他の特徴、ステップ、動作、要素、構成、及び/又はそれらのグループの追加又は存在を排除しない。

【0022】

例示的な実施形態が理想的な実施形態(及び中間構造)の例示的な図面である断面図を参照して、詳細な説明で説明される。このように、例えば、製造技術及び/又は許容誤差のような結果としての図面の形状からの変形が予想されなければならない。つまり、例示的な実施形態は、詳細な説明に示された特定の形状の領域に限定されると解釈されてはならず、例えば、製造からもたらされる形状の偏差を含むべきである。例えば、長方形で示された注入された領域は、一般的には、注入された領域から注入されていない領域へのバイナリ変化というよりは、円形又は曲線のフィーチャー及び/又はエッジでの注入濃度の勾配を有するはずである。同様に、注入によって形成された埋め込みエリアは、埋め込みエリアと注入が行われる表面との間の領域に若干の注入をもたらすことができる。したがって、図面に示された領域は、本質的に概略的であり、その形状は装置の領域の実際の形状を説明するためのものではなく、本発明の範囲を制限しようとするものではない。

【0023】

別の方法で定義されていない限り、本文で使用されるすべての用語(技術的及び科学的用語を含む)は、当業者によって一般的に理解されうる意味を有する。また、一般的に使用される辞典に定義された用語のような用語は、関連技術及び/又は本明細書に関連付けて、その意味と一致する意味を有するものと解釈されるべきであり、本文で定義されていない限り、理想的又は過度に形式的な意味として解釈されてはならない。

【0024】

以下で、例示的な実施形態が添付図面を参照して詳細に説明される。

【0025】

図1は、本発明の実施形態によるシステム100の例示的な実施形態のブロック図であ

10

20

30

40

50

る。図示された実施形態で、システム 100 は、複数のノード又はストレージ装置にわたってデータを格納する分散ストレージシステム 104 を含むことができる。

【0026】

分散ストレージシステムは、大規模な信頼性ストレージ (large-scale reliability storage) を提供する場合、しばしば使用される。多くの場合、これは冗長性 (redundancy) 又はエラー訂正 (例えば、パリティ (parity)) を多数のノード又はストレージ装置にわたって分散させることによって達成される。しかし、ノード又はストレージ装置がオフラインになった場合 (例えば、ネットワークエラー、ハードウェアの故障などに起因する)、データが損傷される可能性があったり、少なくとも冗長性のレベルが減少されたりすることがある。ストレージシステムが分散されるほど、このような状況がより頻繁に発生する。

10

【0027】

このような状況を防止するために、多様な技法 (例えば、ミラーリング、リード・ソロモン符号化) が使用されることができ、本発明の実施形態は、再生成エンコーディングに重点を置く。このような実施形態で、データの欠落している部分 (missing piece) (チャンク (chunk)) は、データの残りの部分に基づいた方式 (formular) を使用して再生成されたり、再構成されたりする。

【0028】

図示された実施形態で、システム 100 は、分散ストレージシステム 104 を管理するように構成された一つ又はそれ以上のホスト装置 102 を含むことができる。ホスト装置 102 は、ストレージシステム 104 を読み書きするコンピューティング装置 (例えば、コンピュータ、サーバ、仮想マシン) を含むことができる。エラー (たとえば、データの欠落しているチャンク) が発生した場合、一般的にホスト装置 102 がエラーを検出し、可能であればエラーを訂正することを担う。

20

【0029】

図示された実施形態で、各データセット 199 は、ホスト装置 102 によって複数のより小さな部分のデータ又はチャンク 198 に分解されたり、又は細分化 (fragmented) されることができる。図示された実施形態で、データ 199 は、チャンク 198 (D1、D2、D3、D4) に分割される。さらに、多様な実施形態で、ホスト装置 102 は、パリティチャンク (P1、P2) (パリティチャンクも、またチャンクであるので、198 の参照番号で表記される。) のような一部の形態の冗長性をデータチャンク 198 に適用することができる。

30

【0030】

本文で、オリジナルデータのチャンク 198 (D1、D2、D3、D4) の個数は、変数 (K 又は k) で説明される。同様に、冗長データのチャンク 198 (P1、P2) の個数は、変数 (R 又は r) で説明される。このようなチャンクの合計数は、 $(K+R)$ 個である。図示された実施形態で、K は 4 であり、R は 2 であり、 $K+R$ は 6 であるが、これは単に例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。

【0031】

図示された実施形態で、ホスト装置 102 は、このようなチャンク 198 (オリジナルと冗長の両方) のそれぞれをストレージシステム 104 のそれぞれのノード又はストレージ装置に格納する。図示された実施形態で、ストレージ装置 114 はチャンク 198 (D1) を格納し、ストレージ装置 114-1 はチャンク 198 (D2) を格納し、ストレージ装置 116 はチャンク 198 (D3) を格納し、ストレージ装置 114-2 はチャンク 198 (D4) を格納し、ストレージ装置 116-1 は、チャンク 198 (P1) を格納し、ストレージ装置 114-3 はチャンク 198 (P2) を格納する。多様な実施形態で、ストレージ装置 (114、116) の個数は、チャンク 198 の個数と同じではなくてよい。

40

【0032】

多様な実施形態で、チャンク 198 が欠落したり、(例えば、ネットワーク又はハードウェアの故障)、又はそうでなければエラーと関連付けられたりする。図示された実施形態で、チャンク 198 (D3) (及びストレージ装置 116) が突然使用できないことを仮

50

定する。ホスト装置 102 は、エラーが検出された場合、チャンク 198 (D3) を再生成したり、そうでなければ、エラーを訂正したりするように試みることができる。

【0033】

このような実施形態で、一つのチャンク (例えば、チャンク 198 (D3)) が故障 (fail) であり、オリジナルデータ 199 に K (例えば、4) 個の全チャンクが存在すれば、故障したチャンク (例えば、チャンク 198 (D3)) を回復するために、少なくとも K (例えば、4) 個のノード又はストレージ装置 (114、116) が、ホスト装置 102 に情報を伝送しなければならない。このような K (例えば、4) 個のチャンクは、(K+R) (例えば、6) のチャンクのいずれから来ることもある。例えば、チャンク 198 (D1、D2、D4、P1) がチャンク 198 (D3) を再生成するために使用されることができる。

10

【0034】

再生成コード (regeneration codes) は、D 個のノード (ただし、一般的には D > K) からの情報を全チャンクサイズの情報よりも少ないサイズで伝送することにより、リペア帯域幅 (repair bandwidth) を減少させる。言い換えると、賢明な方式 (clever formula) の使用により、ホスト装置 102 は、全体のチャンク 198 (D1、D2、D4、P1) を使用せずに、チャンク 198 (D1、D2、D4、P1、P2) のうちただ一部のみを使用することにより、欠落されているチャンク 198 (D3) を再生成することができる。再生成コードは、一般的に、より多くのストレージ装置 (114、116) から情報を取得するものの、非再生成コードが取得するものよりも少ない情報をそれぞれのストレージ装置 (114、116) から獲得する。

20

【0035】

例えば、データの 6 つのチャンクが使用され (K=6)、冗長性の 6 つのチャンクが使用され (R=6、K+R=12)、各チャンクのサイズが 16MB である場合に、標準 RS (Reed-Solomon) エラー訂正エンコーディング方式は、欠落している 16MB のチャンクを訂正するために 6 つ (K) の 16MB のチャンクがホストに伝送されること、すなわち、96MB のデータが伝送されることを必要とする。これと反対に、再生成技法が使用される場合は、12 個すべての (この場合、K+R 又は D) のチャンクの部分がリード (read) されるが、ただ各チャンクの一部 (例えば、2.7MB) のみを使用するため、ホスト装置に伝送される全体のサイズは小さくなる (たとえば、29.7MB)。

【0036】

再生成コードは、ストレージ及び帯域幅のトレードオフを伴う。多様な実施形態で、一般的に再生成コードの 2 つのクラス又はグループが存在する。ストレージのオーバーヘッドが最小である場合、それらは MSR (Minimum Storage Regeneration) コードと呼ばれる。追加されるストレージオーバーヘッドについてリペア帯域幅が最小である場合、それらは MBR (Minimum Bandwidth Regeneration) コードと呼ばれる。このような広いカテゴリーで、多様な特定の技法又は方式が再生成コードを実行するために使用されることができる。前述された内容は、単純な一部の例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。

30

【0037】

図 1 を再び参照すると、図示された実施形態で、ストレージシステム 104 は、複数のストレージ装置 (114、116) を含むことができる。各ストレージ装置 (114、116) は、チャンクにて、又はそれ以外にて、データを格納するように構成されることができる。図示された実施形態で、ストレージ装置 114 は、ハード・ディスク・ドライブ、ソリッド・ステート・ドライブ、又は揮発性メモリのような相対的に従来式のストレージ装置であり得る。

40

【0038】

しかし、図示された実施形態で、ストレージシステム 104 は、再生成コード認識 (RCA: regeneration-code-aware) ストレージ装置 116 を含むことができる。斯かる実施形態で、従来の又は非 RCA ストレージ装置 114 とは異なり、RCA ストレージ装置 116 は、データ再生成コードの演算を支援するように構成されたり、又はそのようにする構成

50

要素を含んだりすることができる。以下で、より詳細に説明されるように、ホスト装置 102 は、データ再生コードの演算の一部を RCA ストレージ装置 116 に動的にオフロード (offload) させることができる。多様な実施形態で、これは、ホスト装置 102 とストレージシステム 104 との間で送受信されるメッセージのサイズ、ホスト装置 102 とストレージシステム 104 との間で伝送されるデータのサイズ、及び/又はホスト装置 102 における演算負荷を減少させることができる。前述された内容は、単に一部の例示的な実施形態であり、本発明の範囲がこれに限定されないことが理解されるだろう。

【0039】

多様な実施形態で、RCA ストレージ装置 116 は、ホスト装置 102 が、最新の又は意図した再生コード方式又は技法でそれらをアップデートできるようにプログラムされることができ、10
斯かる実施形態で、RCA ストレージ装置 116 は、多様な再生技法を格納することができ、再生技法のうち一つは、ホスト装置 102 によって動的に又は半固定的に (semi-statically) 選択されるようにできる。斯かる実施形態で、ホスト装置 102 は、現在どのような再生技法を使用できるかを選択することができる。

【0040】

多様な実施形態で、ストレージシステム 104 は、分散されることができ、10
斯かる実施形態で、ストレージ装置 (114、116) は、互いに物理的に離隔されることができ、ネットワークプロトコルを通じて通信することができる。他の実施形態で、ストレージ装置 (114、116) は、相対的にローカル (例えば、サーバファーム又は同じ建物内) 化されることができ、相変わらずネットワークプロトコルを通じて通信することができ、20
また、別の実施形態で、ストレージシステム 104 は、分散されていなくてもよい。斯かる実施形態で、本発明の実施形態は、ネットワークプロトコル (例えば、USB、SATA) を使用していないローカル装置 (例えば、同一マシン) のために使用されることができ、前述された内容は、単に一部の例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。

【0041】

多様な実施形態で、再生コード認識 (RCA: regeneration-code-aware) のストレージ装置 116 は、再生コードの他のタイプ又はバージョンを演算する能力を含むことができる。30
斯かる実施形態で、望みの再生コードのタイプ又はバージョンは、ホスト装置 102 によって動的に選択されることができ、一部の実施形態で、RCA ストレージ装置 116 は、データをより小さなブロック又はパケットに分割し、イレイジャークード又はそれらの部分を演算し、他の故障であるチャンクのリペアのためにデータチャンク (複数のチャンク) を処理することなどを行うことができる。

【0042】

多様な実施形態で、通信プロトコルは、再生コード又は技法を使用してデータの信頼性を具現するために、ホスト装置 102 と RCA ストレージ装置 116 との間に存在できる。40
斯かる実施形態で、プロトコルは、再生技法を選択すること、入力を渡すこと、意図された技法の動作を指示すること、及び出力を取り出すことを可能にする。一部の実施形態で、プロトコルは、RCA と非 RCA の両方のストレージ装置 (114/116) を含む組み合わせられた環境でホストが動作する際のホスト挙動 (host behavior) 及びそれら両方と相互作用する方法を定めることができる。多様な実施形態で、ホストシステム 102 は、プロトコルを使用して RCA ストレージ装置 116 を設定し、ユーザーのデータをエンコーディングし/リード (read) し/ライト (write) し、データリペアにおける演算をオフロードし、データトラフィックを減少させ、演算を加速し、RCA ストレージ装置の機能を使用してオリジナルデータを再構成することができる。

【0043】

図 2a は、本発明の実施形態によるシステム 201 の例示的な実施形態のブロック図である。50
図示された実施形態で、システム 201 は、第 1 タイプ (Type 1) の再生コードを演算するホスト装置 210 とストレージ装置 212 との間の相互作用 (interaction) を示している。多様な実施形態で、システム 201 は、従来の又は非 RCA ストレージ装置に

使用されることができ、RCA機能が使用されないなら、RCAストレージ装置でも使用されることができる。

【0044】

図示された実施形態で、システム201は、ホスト装置210及びストレージ装置212を含むことができる。このような実施形態で、ホスト装置210は、コマンドを実行して演算を遂行するプロセッサ232と、データ又はデータの一部を少なくとも一時的に格納するメモリ234と、ストレージ装置212又はより一般的にストレージシステム（図示せず）と通信するインターフェース236と、を含むことができる。斯かる実施形態で、ストレージ装置212は、データを格納するように構成されたメモリ224を含むことができる。多様な実施形態で、このようなメモリ224は、不揮発性又は揮発性であり得る。

10

【0045】

図示された実施形態で、チャンク214は、ブロック216に再び分割される。斯かる実施形態で、ホスト装置210は、ストレージ装置212に格納された一つ又はそれ以上のチャンク214から（そして他のストレージ装置に格納されたK-1個のチャンクから）ブロック216を獲得し、再生成コード218（R1）を演算することができる。

【0046】

このような再生成コードの技法（Type 1）において、ブロック216は、より小さなパケット（図示せず）で構成されることができる。各ノード又はストレージ装置212に対し、ホスト装置210は、多様なパケットを使用して、パリティパケット又は再生成コード218を演算する。欠落したり、又はエラーが起きたりしたチャンク（missing or errored chunk）を再構成するために、各ストレージ装置のそれぞれの再生成コード218が使用される。一般的に、第1タイプ（Type 1）の再生成コード技法について、演算は線形であり、故障したチャンクに依存する。再び伝送されるデータのサイズは、サブパケット化（sub-packetization）レベル及び機能に依存する。

20

【0047】

図示された実施形態で、ホスト装置210がエラーを検出した場合、ホスト210は、データリード要請又はコマンド242Aをストレージ装置212に伝送することができる。データリードコマンド242Aは、どのチャンク214がリードされるかに関する情報（たとえば、チャンク214C）を含むことができる。次いで、ストレージ装置212は、データリード応答又はメッセージ244Aを通して、意図されたチャンク214をホスト装置210に伝送する。多様な実施形態で、これは、ホストとストレージ装置212との間の従来のプロトコル（例えば、SATA）を使用して、すべて行われることができる。

30

【0048】

意図されたチャンク214Cが受信されると、ホスト装置210は、インターフェース236によってチャンク214C又はブロック216をメモリ234に格納することができる。プロセッサ232は、次いで、意図された再生成コードの技法287を遂行することができる。再生成コードの技法287は、単純な加算又はブーリアン型のXOR演算（Boolean XORing）に図示されるが、これは単なる一部の例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。前述されたように、多様な実施形態で、これはブロック216をもっと小さなパケットに再び分割することを含むことができる。再生成コードの技法287は、再生成コード218（R1）を演算したり、又は生成したりすることができる。再生成コード218（R1）は、他のチャンク又はストレージ装置と関連付けられている再生成コードとともに、エラーの起きたチャンクを再構成したり、又はリペアしたりするために使用されることができる。

40

【0049】

図2bは、本発明の実施形態によるシステム203の例示的な実施形態のブロック図である。図示された実施形態で、システム203は、第1タイプ（Type 1）の再生成コードを演算するホスト装置210とRCAストレージ装置252との間の相互作用を示す。多様な実施形態で、システム203は、RCAストレージ装置についてのみ使用されることがで

50

き、非RCA装置については使用されない。

【0050】

図示された実施形態では、システム201は、ホスト装置210及びRCAストレージ装置252を含むことができる。斯かる実施形態で、ホスト装置210は、コマンドを実行して演算を遂行するプロセッサ232と、データ又はデータの一部を少なくとも一時的に格納するメモリ234と、RCAストレージ装置252又はより一般的にストレージシステム（図示せず）と通信するインターフェース236と、を含むことができる。

【0051】

このような実施形態で、RCAストレージ装置252は、データを格納するように構成されたメモリ224を含むことができる。多様な実施形態で、メモリ224は、不揮発性又は揮発性であり得る。さらに、多様な実施形態で、RCAストレージ装置252は、ホスト装置210（一般的に、ストレージ装置の外部に位置する）によって要請された場合、選択された数のデータブロック216に基づいてデータ再生コード218を演算するように構成されたプロセッサ222を含むことができる。多様な実施形態で、プロセッサ222は、プログラム可能なゲートアレイ（例えば、FPGA）、グラフィックスプロセッサユニット（GPU：graphics processor unit）、汎用プロセッサ（例えば、CPU）、コントローラプロセッサ、又はシステム・オン・チップ（SoC：system-on-chip）を含むことができる。前述された内容は、単に一部の例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。RCAストレージ装置252は、コマンドの複数のセットを格納するように構成されたコードメモリ228を含むことができる。コマンドの各セット229は、他の再生コードの技法を遂行する方法についてのコマンド又は他のデータ再生コードを生成する。多様な実施形態で、コマンドのセット229は、ストレージ装置252に予め構成されることができたり、駆動中に（例えば、ホスト装置210によって）動的に追加/調節できたり、又はそれらの組み合わせの形で具現されたりすることができる。RCAストレージ装置252は、少なくとも、ホスト装置210と通信するように構成された外部インターフェース226を含むことができる。

【0052】

図示された実施形態で、ホスト装置210は、ストレージ装置252が、データ再生コードを内部的に演算できるか、又は一般的にRCAストレージ装置であるかを判定することができる。もしそうであれば、ホスト装置210は、RCAストレージ装置252が意図された再生コードの技法を遂行できるか、又は意図された再生コードの技法を遂行するように（コードメモリ228を介して）プログラムされることができるとかを判定することができる。もしそうでなければ、図2aに示された技法を使用することができる。

【0053】

RCAストレージ装置252が意図された再生コードの技法を遂行することができる場合、ホスト装置210は、リペアのためのリードコマンド242B（Read for Repair command）を発行することができる。多様な実施形態で、リペアのためのリードコマンド242Bは、意図された再生又はリペア技法の指示（indication）、意図されたパケット又はブロックサイズ、意図された再生コード又はリペア技法についてのパラメータ、データ又はチャンクアドレス、故障したチャンクナンバーのいずれか一つ又はそれ以上を含んだり、示したりすることができる。これは単に一部の例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。

【0054】

コマンド242Bに回答して、プロセッサ222は、意図されたブロック216又はチャンク214Cを取り出すことができる。プロセッサ222は、意図された再生又はリペア技法に関連付けられたコマンドのセット229を取り出すことができる。プロセッサ222は、意図された再生技法287を遂行し、データ再生コード218（DRC）（R1）を演算することができる。

【0055】

その後、RCAストレージ装置252は、ホスト装置210へ、インターフェース226

10

20

30

40

50

を介して、データ再生コード 2 1 8 (R1)(メッセージ 2 4 4 B)を伝送することができる。斯かる実施形態で、データ再生コード 2 1 8 (R1)は、図 2 aのメッセージ 2 4 4 Aを通して伝送されるデータよりも小さいサイズを有したり、又はより少ない帯域幅を使用したりすることができる。

【 0 0 5 6 】

図示された実施形態で、メッセージ (2 4 2 B、 2 4 4 B) は、メッセージ (2 4 2 A、 2 4 4 A) に対して使用されたものとは異なるプロトコルを必要とし得る。メッセージ (2 4 2 A、 2 4 4 A) は、従来のストレージ装置のプロトコルによって許容されることができる一方、メッセージ (2 4 2 B、 2 4 4 B) は、追加的であり、他の情報を必要とすることがあり、これにより、新しいメッセージングプロトコル又は少なくとも新しいコマンドを必要とし得る。

10

【 0 0 5 7 】

図示された実施形態で、ホスト装置 2 1 0 は、他のRCAのストレージ装置 (図示せず) によって提供されたり、ホスト装置 2 1 0 自体によって生成されたりした追加的なデータ再生コードとともにデータ再生コード 2 1 8 (R1) を使用して、エラーの起きたデータのチャンクを再生成することができる。

【 0 0 5 8 】

図 2 cは、本発明の実施形態によるシステム 2 0 5 の例示的な実施形態のブロック図である。図示された実施形態で、システム 2 0 5 は、第2タイプ (Type 2) の再生コードを演算するホスト装置 2 1 0 とストレージ装置 2 1 2 との間の相互作用を示している。多様な実施形態で、システム 2 0 5 は、従来の又は非RCAストレージ装置に対して使用されることができる、RCA機能が使用されていない場合、RCAストレージ装置についても使用されることができる。

20

【 0 0 5 9 】

図示された実施形態で、システム 2 0 5 は、ホスト装置 2 1 0 及びストレージ装置 2 1 2 を含むことができる。ホスト装置 2 1 0 とストレージ装置 2 1 2 の両方は、先に説明されて図示された構成要素を含むことができる。

【 0 0 6 0 】

このような再生コードの技法 (Type 2) で、データ再生コードは、より少ないパケット (図示せず) 又はブロック 2 1 6 がリードされるように演算される。しかし、これは、意図されたブロック 2 1 6 又はパケットが事前に完全に知られるが、演算が遂行されるにつれて断片的に要請されることを意味する。斯かるタイプの再生技法は、ネットワーク帯域幅とデータリードの両方を理論上では減少させるが、一つの大きなリード(read)を複数のより小さなリードに変換し、これは性能に良くない。

30

【 0 0 6 1 】

図示された実施形態で、ホスト装置 2 1 0 は、意図された再生技法の一部 2 8 8 を使用して、ブロック (E1) がエラーと関連付けられている場合、ブロック (B1、 B3) (又はそれらのパケット) がエラーの起きたブロック (E1) を修正するのに必要なものであると計算する。斯かる実施形態で、ブロック (B1) が必要であるとホスト装置 2 1 0 が検出すると、ホスト装置 2 1 0 は、データリード要請又はコマンド 2 4 2 C をストレージ装置 2 1 2 に伝送することができる。データリードコマンド 2 4 2 C は、どのブロック 2 1 6 がリードされるか (例えば、ブロック (B1)) についての情報を示すことができる。ストレージ装置 2 1 2 は、次いで、意図されたブロック 2 1 6 (B1) を、データリード応答又はメッセージ 2 4 4 C を通して、ホスト装置 2 1 0 に伝送する。多様な実施形態で、これはホストとストレージ装置との間の従来のプロトコル (例えば、SATA) を使用して、すべて行われることができる。

40

【 0 0 6 2 】

このような実施形態で、ブロック (B3) が必要であるとホスト装置 2 1 0 が検出すると、ホスト装置 2 1 0 は、データリード要請又はコマンド 2 4 6 C をストレージ装置 2 1 2 に伝送することができる。これは、一般的に、ブロック (B1) を要請するものとは別に、

50

第2のデータ要請として遂行される。データリードコマンド246Cは、どのブロック215がリードされるか(例えば、今回はブロック(B3))についての情報を示すことができる。ストレージ装置212は、次いで、データリード応答又はメッセージ248Cを通して意図されたブロック216(B3)をホスト装置210に伝送する。多様な実施形態で、これはホストとストレージ装置との間の従来のプロトコル(例えば、SATA)を使用して、すべて行われることができる。

【0063】

意図されたブロック216が受信されると、インターフェース236によってホスト装置210は、ブロック216をメモリ234に格納することができる。プロセッサ232は、次いで、意図された再生成コードの技法(部分289によって図示される。)を行うことができる。再生成コードの技法(又は部分289)は、再生成コード219(R1)を演算したり、生成したりすることができ、再生成コード219(R1)は、他のチャンク又はストレージ装置と関連付けられている再生成コードとともに、エラーの起きたチャンクを再構成したり、リペアしたりするのに使用されることができる。

【0064】

図2dは、本発明の実施形態によるシステム207の例示的な実施形態のブロック図である。図示された実施形態で、システム207は、第2タイプ(Type 2)の再生成コードを演算するホスト210とRCAストレージ装置252との間の相互作用を示している。多様な実施形態で、システム207は、RCAストレージ装置についてのみ使用されることができ、非RCAストレージ装置については、使用されない。

【0065】

図示された実施形態で、システム207は、ホスト装置210及びストレージ装置252を含むことができる。ホスト装置210とストレージ装置252の両方は、先に説明されて図示された構成要素を含むことができる。

【0066】

図示された実施形態で、ホスト装置210は、ストレージ装置252が、データ再生成コードを内部的に演算することができるか、又は一般的にRCAストレージ装置であるかを判定することができる。もしそうであれば、ホスト装置210は、RCAストレージ装置252が意図された再生成コードの技法を遂行できるか、又は意図された再生成コードの技法を遂行するように(コードメモリ228を介して)プログラムされることができるかを判定することができる。もしそうでなければ、図2Cに示された技法を使用することができる。

【0067】

RCAストレージ装置252が意図された再生成コードの技法を遂行することができる場合、ホスト装置210は、リペアのためのリードコマンド242D(Read for Repair command)を発行することができる。多様な実施形態で、リペアのためのリードコマンド242Dは、意図された再生成又はリペア技法の指示(indication)、意図されたバケット又はブロックサイズ、意図された再生成又はリペア技法のためのパラメータ、データ又はチャンクアドレス、故障したチャンクナンバー(例えば、ブロック(E1))のいずれか1つ又はそれ以上を含んだり、示したりすることができる。前述された内容は、単に一部の例示的な実施形態であり、本発明の範囲がこれに限定されないことが理解されるだろう。

【0068】

コマンド242Dに応答して、プロセッサ222は、意図された再生成又はリペア技法に関連付けられたコマンドのセット229を取り出すことができる。プロセッサ222は、意図された再生成技法又はその一部分288を遂行することができる。斯かる実施形態で、プロセッサ222は、意図されたブロックがB1及びB3であると演算することができる。斯かる実施形態で、斯かるブロック(B1、B3)は、RCAストレージ装置252によって演算されたデータ再生成コードに含められることができる。斯かる実施形態で、斯かるブロックは、リペアのためのリードコマンド242Dの応答の一部と見なされうる。

【0069】

10

20

30

40

50

RCAストレージ装置 252 は、次いで、インターフェース 226 を介して意図されたブロック (B1、B3) (例えば、メッセージ 244D) をホスト装置 210 に伝送することができる。このような実施形態で、データ再生コード又は意図されたブロック (B1、B3) は、より小さなサイズを有するか、より少ない帯域幅を使用するか、又は少なくともより小さなメッセージを含むかであり得る、これにより図 2c のメッセージ (244C、248C) を通して伝送されたデータよりもオーバーヘッドが少なくなり得る。

【0070】

図示された実施形態で、メッセージ (242D、244D) は、メッセージ (242C、244C、246C、248C) のために使用されているものとは異なるプロトコルを必要とし得る。メッセージ (242C、244C、246C、248C) は、従来のストレージ装置のプロトコルによって許容されることができ、メッセージ (242D、244D) は、他の追加的な情報を必要とすることがあり、これにより、新しいメッセージプロトコル又は少なくとも新しいコマンドが要求され得る。

10

【0071】

図示された実施形態で、ホスト装置 210 は、次いで、データ再生コード又はブロック (B1、B3) を、他の RCA ストレージ装置 (図示せず) によって提供されたり、又はホスト装置 210 自体によって生成されたりした追加の再生コード又はデータとともに使用して、エラーの起きたデータ (E1) を再生成することができる。

【0072】

図 3 は本発明の実施形態による技法 300 の例示的な実施形態のフローチャートである。多様な実施形態で、技法 300 は、図 1、図 2a、図 2b、図 2c 及び図 2d のようなシステムによって使用されたり、又は生成されたりすることができる。前述された内容は、単に一部の例示的な実施形態であり、本発明がこれに限定されないことが理解されるだろう。本発明は、技法 300 によって図示された多数の動作の順序に限定されないことが理解されるだろう。

20

【0073】

図示された実施形態で、説明の便宜のために、技法 300 は、ストレージシステムのすべての装置が RCA ストレージ装置又は非 RCA ストレージ装置のいずれか一方である実施形態 (つまり、同種のストレージシステム (homogeneous storage system)) を示す。組み合わせ又は異種のストレージシステム (heterogeneous storage system) については、本発明が属する技術分野における当業者は、単純化された技法 300 が個別的なストレージ装置に適用されるようにどのように拡張されるかをよく理解するはずである。

30

【0074】

一実施形態で、ブロック 302 は、データのチャンクと関連付けられたエラーが検出されうることを示している。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作 (複数の動作) は、図 1、図 2a、図 2b、図 2c、及び図 2d のシステム又は装置によって遂行されることができ。

【0075】

一実施形態で、ブロック 304 は、前述されたように、データ再生コード (DRC) がホスト装置によって演算されるのか、又はそれぞれの RCA ストレージ装置によって演算されるのかが判定されることを示す。多様な実施形態で、このようなブロックによって図示された一つ又はそれ以上の動作 (複数の動作) は、前述されたように、図 1、図 2a、図 2b、図 2c 及び図 2d のシステム又は装置によって遂行されることができ。

40

【0076】

一実施形態で、ブロック 306 は、DRC がホストによって、より一般的な方式で演算される場合に、DRC を演算するのに十分な既存のデータが存在するのかが判定されうることを示す。斯かる実施形態で、これは K 個のチャンクが (K+R) 個のデータチャンクから (out of the (K+R) data chunks) 利用可能であるのかを判定することを含むことができる。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作 (複数の動作) は、前述されたように、図 1、図 2a、図 2b、図 2c 及び図 2d のシステム又は装置

50

によって遂行されることができる。

【 0 0 7 7 】

一実施形態で、ブロック 3 9 9 は、DRCを演算するのに十分なエラーフリーのチャンクが存在しない場合に、エラーであるデータのチャンクの再生成を超えて他の形態のエラー処理が発生できることを示す。多様な実施形態で、これは単にデータがエラーを有したり、又は利用可能ではなかったりすることを報告するものである。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c及び図 2 dのシステム又は装置によって遂行されることができる。

【 0 0 7 8 】

一実施形態で、ブロック 3 0 8 は、必要とされる個数のチャンク（例えば、K個のチャンク）が、多様な（例えば、K+R個）ストレージ装置からリードされうることを示す。多様な実施形態で、このようなブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c及び図 2 dのシステム又は装置によって遂行されることができる。

【 0 0 7 9 】

一実施形態で、ブロック 3 1 0 は、前述されたように、ホスト装置がエラーフリーのチャンク（例えば、K個のチャンク）を使用して、エラーの起きたチャンクを再構成したり、又は再生成できたりすることを示す。多様な実施形態で、斯かるブロックによって図示された一つ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c及び図 2 dのシステム又は装置によって遂行されることができる。

【 0 0 8 0 】

一実施形態で、ブロック 3 5 0 は、前述されたように、エラーフリーのチャンク（例えば、Dのチャンク）がDRCを演算するのに十分に存在するのかが判定されうることを示す。もしそうでなければ、多様な実施形態で、技法 3 0 0 は、ブロック 3 0 6 から始まる非RCA装置のパスを試みることに頼られる。それ以外は、技法 3 0 0 は、ブロック 3 5 2 に続くことができる。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c及び図 2 dのシステム又は装置によって遂行されることができる。

【 0 0 8 1 】

一実施形態で、ブロック 3 5 2 は、前述されたように、リペアのためのリードコマンド（read for repair command）が全体（例えば、K+R個）のストレージ装置のうち、必要な個数（例えば、D個）に発行されうることを示している。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c及び図 2 dのシステム又は装置によって遂行されることができる。

【 0 0 8 2 】

一実施形態で、ブロック 3 5 4 は、前述されたように、複数のバージョン又はタイプのDRC技法のうちのどれが使用されるのかが判定されうることを示している。図示された実施形態で、DRC技法のバージョン及びタイプは、前述されたように、第1タイプ（Type 1）及び第2タイプ（Type 2）に一般化される。しかし、このようなタイプは、単に一部の例示的な実施形態であり、本発明がこれに限定されるものではなく、また、広い範囲のタイプにおいて、前述されたように、多くのサブタイプが存在できることが理解されるだろう。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c及び図 2 dのシステム又は装置によって遂行されることができる。

【 0 0 8 3 】

一実施形態で、ブロック 3 5 6 は、第1タイプのDRC手法（Type 1 DRC technique）が選択された場合、前述されたように、リペア機能（repair function）がRCAストレージ装置内でチャンクに適用されうることを示す。多様な実施形態で、斯かるブロックに

10

20

30

40

50

よって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c 及び図 2 d のシステム又は装置によって遂行されることができる。

【 0 0 8 4 】

一実施形態で、ブロック 3 5 8 は、第 2 タイプの DRC 手法（Type 2 DRC technique）が選択された場合、前述されたように、リペアのために要求されるブロック（又はパケットのような他のサブ部分）が演算されうることを示す。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c 及び図 2 d のシステム又は装置によって遂行されることができる。

【 0 0 8 5 】

一実施形態で、ブロック 3 6 0 は、DRC 又は要求されるブロックが演算されると、前述されたように、DRC 又はブロックがホスト装置に伝送されうることを示す。多様な実施形態で、これは、前述されたように、非 RCA パスよりも小さいサイズのデータ又はより小さい個数のメッセージを含むことができる。多様な実施形態で、斯かるブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c 及び図 2 d のシステム又は装置によって遂行されることができる。

【 0 0 8 6 】

一実施形態で、ブロック 3 6 2 は、前述されたように、ホスト装置が DRC 又は返還されたブロックを使用して、エラーの起きたチャンクを再構成したり、又は再生成できたりすることを示す。多様な実施形態で、このようなブロックによって図示された一つ又はそれ以上の動作（複数の動作）は、前述されたように、図 1、図 2 a、図 2 b、図 2 c 及び図 2 d のシステム又は装置によって遂行されることができる。

【 0 0 8 7 】

多様な実施形態で、ホスト装置は、データ再生成コードを演算する RCA ストレージ装置（又はそれに含まれたプロセッサ）の機能をオン又はオフさせるコマンドを RCA ストレージ装置に提供することができる。RCA ストレージ装置（又はそれに含まれたプロセッサ）は、コマンドに応答して、図 2 b 又は図 2 d に図示された RCA ストレージ装置のように動作したり、又は図 2 a 又は図 2 c に図示された従来の又は非 RCA ストレージ装置のように動作できたりする。

【 0 0 8 8 】

図 4 は、本発明の技術的思想に基づいて形成された半導体装置を含むことができる情報処理システム 4 0 0 の概念的なブロック図である。

【 0 0 8 9 】

図 4 を参照すると、情報処理システム 4 0 0 は、本発明の技術的思想に基づいて構成された一つ又はそれ以上の装置を含むことができる。他の実施形態で、情報処理システム 4 0 0 は、本発明の技術的思想に基づいた一つ又はそれ以上の技法を使用したり、又は実行できたりする。

【 0 0 9 0 】

多様な実施形態で、情報処理システム 4 0 0 は、例えば、ラップトップ（laptop）、デスクトップ、ワークステーション、サーバ、ブレードサーバ、パーソナルデジタルアシスタント（PDA：personal digital assistant）、スマートフォン、タブレット、及び他の適切なコンピュータのようなコンピューティング装置又は仮想マシン若しくはそれらの仮想コンピューティング装置を含むことができる。多様な実施形態で、情報処理システム 4 0 0 は、ユーザー（図示せず）によって使用されることができる。

【 0 0 9 1 】

本発明による情報処理システム 4 0 0 は、中央処理ユニット（CPU：central processing unit）、ロジック又はプロセッサ 4 1 0 をさらに含むことができる。一部の実施形態で、プロセッサ 4 1 0 は、一つ又はそれ以上の機能ユニットブロック（FUBs）若しくは組み合わせロジックブロック（CLBs）4 1 5 を含むことができる。このような実施形態で、組み合わせロジックブロックは、多様なブーリアン方式のロジック動作（例えば、NAND

10

20

30

40

50

、NOR、NOT、XOR）、安定化ロジックデバイス（例えば、フリップフロップ、ラッチ）、他のロジックデバイス、又はそれらの組み合わせを含むことができる。このような組み合わせロジックの動作は、単純な又は複雑な方式で入力信号を処理して、意図された結果を達成するように構成されることができる。同期式組み合わせロジックの動作の一部の例示的な実施形態が説明されたが、本発明がこれに限定されるものではなく、非同期式動作又はそれらの組み合わせを包含できることが理解されるだろう。一実施形態で、組み合わせロジックの動作は、複数のCMOS（complementary metal oxide semiconductor）トランジスタを含むことができる。多様な実施形態で、斯かるCMOSTランジスタは、ロジック動作を遂行するゲートへと構成されることができるが、ただし、本発明の範囲に属する他の技術が使用されうることが理解される。

10

【0092】

本発明による情報処理システム400は、揮発性メモリ420（例えば、ランダムアクセスメモリ（RAM：random access memory））をさらに含むことができる。本発明による情報処理システム400は、不揮発性メモリ430（例えば、ハードドライブ、光メモリ、NAND又はフラッシュメモリ）をさらに含むことができる。一部の実施形態で、揮発性メモリ420、不揮発性メモリ430、又はそれらの組み合わせ若しくは一部は、「ストレージ媒体（storage medium）」と称されることができる。多様な実施形態で、揮発性メモリ420及び/又は不揮発性メモリ430は、半永久的な又は実質的に永久的な形態でデータを格納するように構成されることができる。

【0093】

多様な実施形態で、情報処理システム400は、情報処理システム400が通信ネットワークの一部になるように、又は通信ネットワークを介して通信するように構成された一つ又はそれ以上のネットワークインターフェース440を含むことができる。Wi-Fiプロトコルの例示は、IEEE（Institute of Electrical and Electronics Engineers）802.11g、IEEE 802.11nを包含できるが、これに限定されない。セルラープロトコルの例示は、IEEE 802.16m（別名、Wireless-MAN（Metropolitan Area Network）Advanced）、LTE（Long Term Evolution）Advanced、Enhanced Data rates for GSM（登録商標（Global System for Mobile Communications））Evolution（EDGE）、Evolved High-Speed Packet Access（HSPA+）を包含できるが、これに限定されない。有線プロトコル例示は、IEEE 802.3（別名、Ethernet）、ファイバチャネル（Fibre Channel）、電力線通信（Power Line communication）（例えば、ホームプラグ、IEEE 1901）を包含できるが、これに限定されない。前述された内容は、単に一部の例示的な実施形態であり、本発明はこれに限定されないことが理解されるだろう。

20

30

【0094】

本発明による情報処理システム400は、ユーザーインターフェースユニット450（例えば、ディスプレイアダプター、ハプティックインターフェース、ヒューマンインターフェースデバイス）をさらに含むことができる。多様な実施形態で、このようなユーザーインターフェースユニット450は、ユーザーから入力を受信したり、又はユーザーに出力を提供したりするように構成されることができる。他の種類の装置がユーザーとの相互作用を提供するために使用されることができ、例えば、ユーザに提供されるフィードバックは、感覚フィードバック（sensory feedback）（例えば、視覚フィードバック、聴覚フィードバック又は触覚フィードバック）の形態であり得るし、ユーザからの入力は、音、音声又は触覚入力の形態で受信されることができる。

40

【0095】

多様な実施形態で、情報処理システム400は、1つ又はそれ以上の他の装置又はハードウェア構成460（例えば、ディスプレイ又はモニター、キーボード、マウス、カメラ、指紋リーダー、ビデオプロセッサ）を含むことができる。前述された内容は、単に一部の例示的な実施形態であり、本発明はこれに限定されないことが理解されるだろう。

【0096】

本発明による情報処理システム400は、1つ又はそれ以上のシステムバス405をさ

50

らに含むことができる。斯かる実施形態で、システムバス405は、プロセッサ410、揮発性メモリ420、不揮発性メモリ430、ネットワークインターフェース440、ユーザインターフェースユニット450、及び1つ又はそれ以上のハードウェア構成460を通信的に連結するように構成されることができる。プロセッサ410によって処理されたデータ又は不揮発性メモリ430の外部から入力されたデータは、不揮発性メモリ430又は揮発性メモリ420のいずれか一つに格納されることができる。

【0097】

多様な実施形態で、情報処理システム400は、1つ又はそれ以上のソフトウェア構成470を含んだり、又は実行したりすることができる。一部の実施形態で、ソフトウェア構成470は、オペレーティングシステム(OS: operating system)及び/又はアプリケーションを含むことができる。一部の実施形態で、OSは一つ又はそれ以上のサービスをアプリケーションとして提供し、アプリケーションと情報処理システム400の多様なハードウェア構成(例えば、プロセッサ410、ネットワークインターフェース440)との間の仲裁者(intermediary)として管理したり、又は動作したりすることができる。このような実施形態で、情報処理システム400は、地域的に(例えば、不揮発性メモリ430内に)インストールされることができ、プロセッサ410によって直接的に実行され、OSと直接相互作用するように構成された一つ又はそれ以上のネイティブアプリケーションを含むことができる。このような実施形態で、ネイティブアプリケーションは、プリコンパイルされたマシン実行可能なコード(pre-compiled machine executable code)を含むことができる。一部の実施形態で、ネイティブアプリケーションは、ソース又はオブジェクトコードをプロセッサ410によって実行される実行可能なコードに変換するように構成されたVM(virtual execution machine)(例えば、the Java Virtual Machine、the Microsoft Common Language Runtime)又はスクリプトインタプリタ(script interpreter)(例えば、csh(C shell)、AppleScript、AutoHotkey)を含むことができる。

【0098】

前述した半導体装置は、多様なパッケージング技法を使用してカプセル化されることができる。例えば、本発明の技術的思想に基づいて構成された半導体装置は、POP(package on package)技法、BGA(ball grid array)技法、CSP(chip scale package)技法、PLCC(plastic leaded chip carrier)技法、PDIP(plastic dual in-line package)技法、ダイインワッフルパック技法(die in waaffle pack technique)、ダイインウェハー形成技法(die in wafer form technique)、COB(chip on board)技法、CERDIP(ceramic dual in-line package)技法、PMQFP(plastic metric quad flat package)技法、PQFP(plastic quad flat package)技法、SOLIC(small outline package)技法、SSOP(shrink small outline package)技法、TSOP(thin small outline package)技法、TQFP(thin quad flat package)技法、SIP(system in package)技法、MCP(multi-chip package)技法、WFP(wafer-level fabricated package)技法、WSP(wafer-level processed stack package)技法、当業者によってよく知られている他の技法のいずれか一つを使用してカプセル化されることができる。

【0099】

方法のステップは、コンピュータプログラムを実行して入力データに対して動作したり、又は出力を生成したりすることにより、機能を遂行する一つ又はそれ以上のプログラム可能なプロセッサによって遂行されることができる。方法のステップは、特定の目的のためのロジック回路、例えば、FPGA(field programmable gate array)又はASIC(application-specific integrated circuit)によって遂行されることができ、装置は特定の目的のためのロジック回路、例えば、FPGA(field programmable gate array)又はASIC(application-specific integrated circuit)によって具現されることができる。

【0100】

多様な実施形態で、コンピュータ読み取り可能な媒体は、コマンドを含むことができ、コマンドが実行される場合、装置は、方法のステップの少なくとも一部を遂行することができる。一部の実施形態で、コンピュータ読み取り可能な媒体は、磁気媒体 (magnetic medium)、光媒体 (optical medium)、他の媒体、又はそれらの組み合わせ (例えば、CD-ROM、ハードドライブ、リードオンリーメモリ、フラッシュドライブ) に含まれることができる。このような実施形態で、コンピュータ読み取り可能な媒体は、明確かつ非一時的に具体化された製造品であり得る。

【 0 1 0 1 】

本発明の技術的思想が例示的な実施形態を参照して説明されたが、当業者は、本発明の思想及び範囲から逸脱せずに多様な変形及び変化を行うことができる。したがって、前述した実施形態は、ただ説明のためのものであり、制限されないことが理解されるだろう。本発明の範囲は以下の特許請求の範囲及びその均等物の最も広く許容される解釈によって決定されるべきで、前記の説明によって制限されたり、又は限定されたりしてはならない。したがって、添付された特許請求の範囲は、実施形態の範囲内でこのような変形と修正の両方を含むものと意図されることが理解されるだろう。

【 符号の説明 】

【 0 1 0 2 】

1 0 0、2 0 1、2 0 3、2 0 5、2 0 7 システム
1 0 2、2 1 0 ホスト装置
1 0 4 ストレージシステム
1 1 4、2 1 2 (非RCA) ストレージ装置
1 1 6、2 5 2 RCA ストレージ装置
1 9 8、2 1 4 チャンク
1 9 9 データ
2 1 6 ブロック
2 1 8、2 1 9 データ再生成コード
2 2 2 プロセッサ
2 2 4 メモリ 2 2 6 外部インターフェース
2 2 8 コードメモリ

10

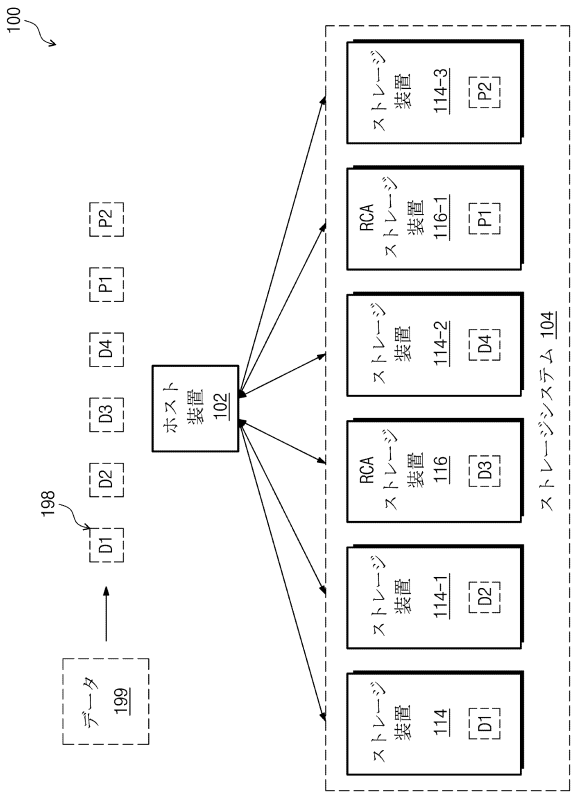
20

30

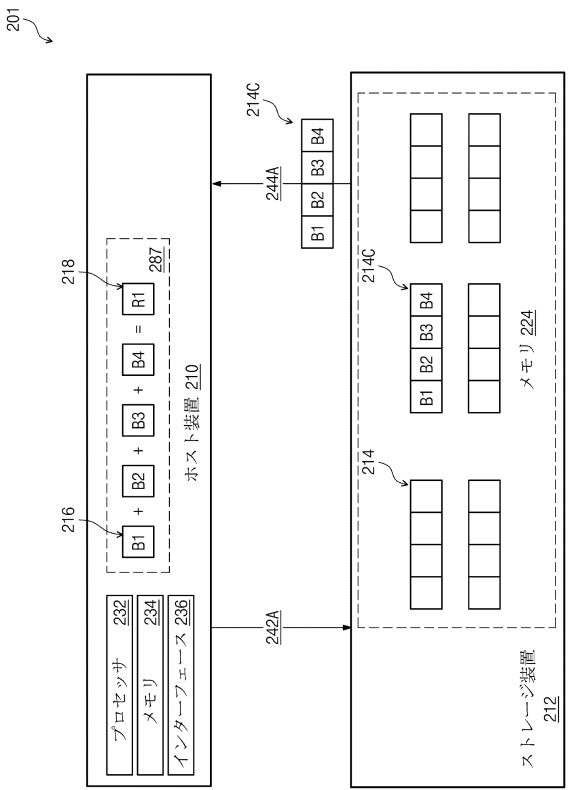
40

50

【図面】
【図 1】



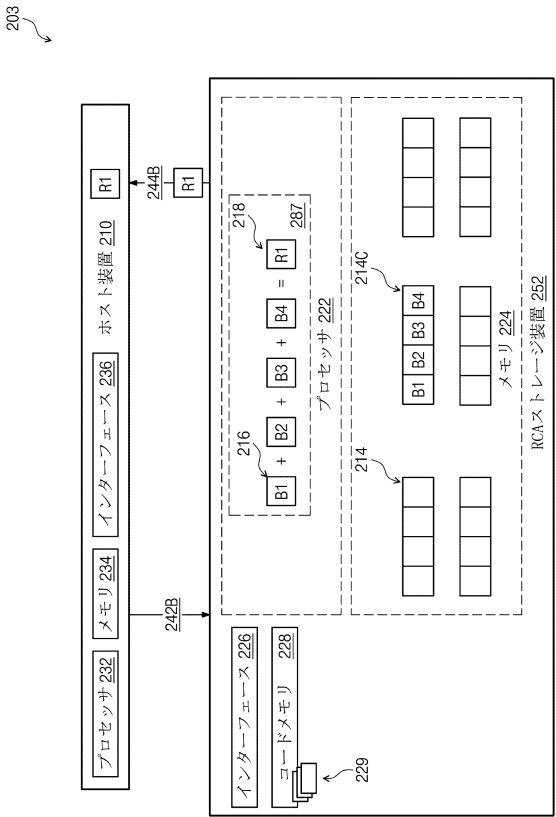
【図 2 a】



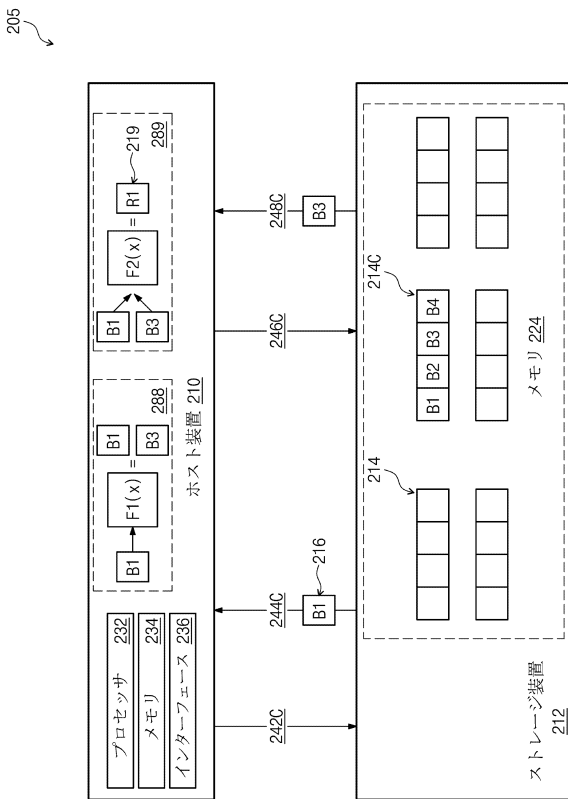
10

20

【図 2 b】



【図 2 c】

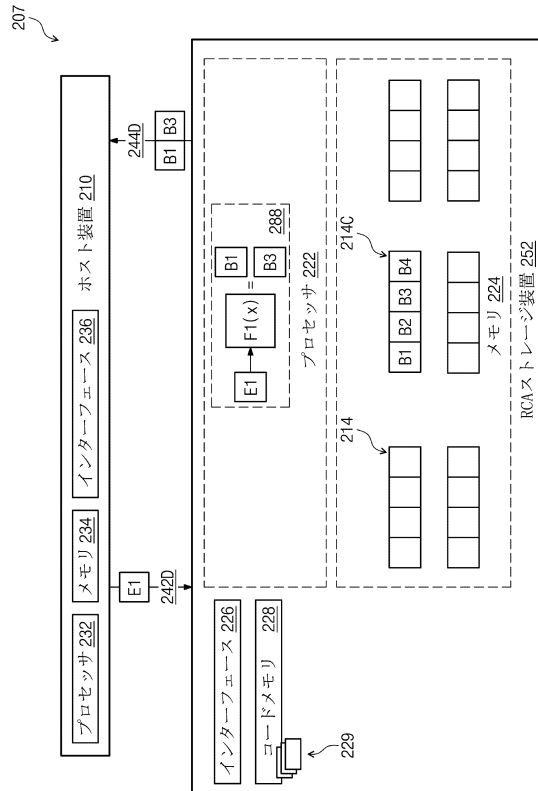


30

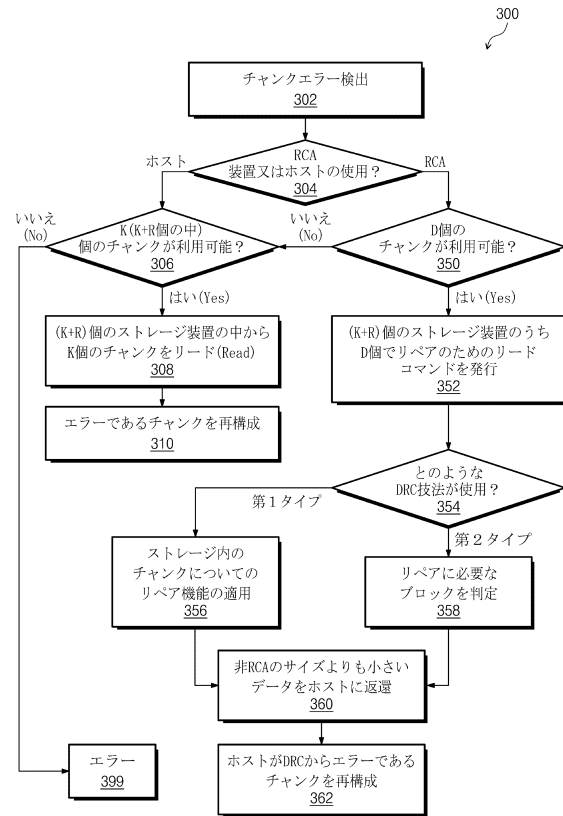
40

50

【 図 2 d 】



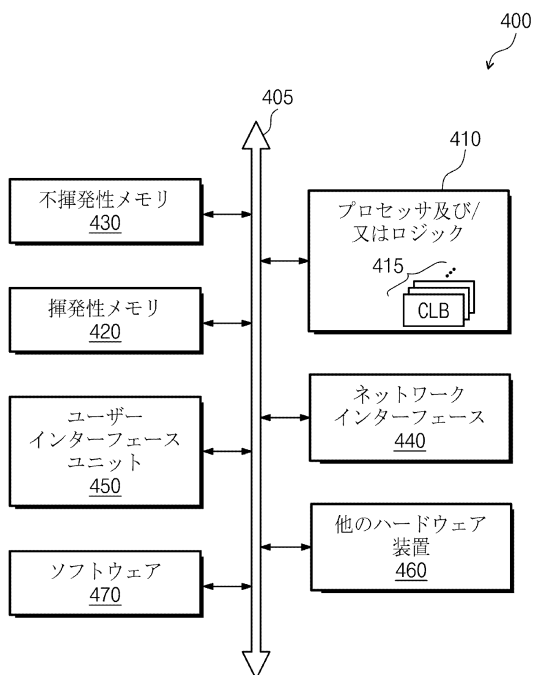
【 図 3 】



10

20

【 図 4 】



30

40

フロントページの続き

弁理士 大貫 進介
(72)発明者 レッカ ピチュマニ
アメリカ合衆国 バージニア州 22033 フェアファックス ウェイザーン プレイス 3871
(72)発明者 奇 亮 そく
アメリカ合衆国 カリフォルニア州 94303 パロアルト アルテールウォーク 873
審査官 田中 幸雄
(56)参考文献 国際公開第2017/065628(WO, A1)
米国特許出願公開第2017/0161148(US, A1)
(58)調査した分野 (Int.Cl., DB名)
G06F 11/10
G06F 3/06
G06F 13/10