

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4097436号
(P4097436)

(45) 発行日 平成20年6月11日(2008.6.11)

(24) 登録日 平成20年3月21日(2008.3.21)

(51) Int.Cl.

F I

G 0 6 F 11/20 (2006.01)

G 0 6 F 11/20 3 1 0 K

請求項の数 8 (全 23 頁)

(21) 出願番号	特願2002-45800 (P2002-45800)	(73) 特許権者	398038580
(22) 出願日	平成14年2月22日(2002.2.22)		ヒューレット・パカード・カンパニー
(65) 公開番号	特開2002-328815 (P2002-328815A)		HEWLETT-PACKARD COMPANY
(43) 公開日	平成14年11月15日(2002.11.15)		アメリカ合衆国カリフォルニア州パロアルト
審査請求日	平成17年1月19日(2005.1.19)		ハノーバー・ストリート 3000
(31) 優先権主張番号	09/810, 103	(74) 代理人	100081721
(32) 優先日	平成13年3月15日(2001.3.15)		弁理士 岡田 次生
(33) 優先権主張国	米国 (US)	(74) 代理人	100105393
			弁理士 伏見 直哉
		(74) 代理人	100111969
			弁理士 平野 ゆかり

最終頁に続く

(54) 【発明の名称】 冗長コントローラシステムからコントローラをオンライン除去する方法

(57) 【特許請求の範囲】

【請求項 1】

第 1 のメモリを有する第 1 のコントローラおよび第 2 のメモリを有する第 2 のコントローラを有する冗長コントローラシステムから、オンラインであるコントローラを除去するための方法であって、

早期検出信号に応じて、前記第 1 のコントローラが、前記冗長コントローラシステムから該第 1 のコントローラの除去が行われていることを検出するステップであって、該早期検出信号は、前記第 1 のコントローラの前記システムからの除去が行われている間にわたって発せられる、ステップと、

前記第 1 のコントローラおよび前記第 2 のコントローラの両方が、それぞれの前記第 1 および第 2 のメモリに対する未処理のメモリアクセスを完了させることを含むシャットダウンシーケンスを実行するステップと、

前記第 1 のコントローラが、前記第 1 のメモリを自己リフレッシュモードにすると共に、前記第 2 のコントローラが、前記第 2 のメモリを自己リフレッシュモードにするステップと、

前記第 1 のコントローラが、オフラインに留まって前記冗長コントローラシステムからの前記第 1 のコントローラの除去が完了するのを待つと共に、前記第 2 のコントローラが、オンラインになる処理を開始するステップと、

を含む方法。

【請求項 2】

10

20

前記早期検出信号に応じて前記第 1 のコントローラの第 1 のプロセッサに割り込みがかけられると共に、前記第 2 のコントローラの第 2 のプロセッサに割り込みがかけられることによって前記シャットダウンシーケンスは実行される、

請求項 1 に記載の方法。

【請求項 3】

前記第 1 のプロセッサに割り込みがかけられることにより、前記第 1 のプロセッサによる未処理のプロセッサタスクが終了されると共に、前記第 2 のプロセッサに割り込みがかけられることにより、前記第 2 のプロセッサによる未処理のプロセッサタスクが終了される、

請求項 2 に記載の方法。

10

【請求項 4】

前記第 1 のメモリを前記自己リフレッシュモードにすることは、前記第 1 のメモリを、メモリの内容を維持するためのメモリコントローラからの外部リフレッシュサイクルを必要としないモードにすることを含む、

請求項 1 に記載の方法。

【請求項 5】

前記第 2 のメモリを自己リフレッシュモードにすることは、前記第 2 のメモリを、メモリの内容を維持するためのメモリコントローラからの外部リフレッシュサイクルを必要としないモードにすることを含む、請求項 1 に記載の方法。

【請求項 6】

20

前記第 2 のメモリの自己リフレッシュモードを終了させた後で、前記第 2 のコントローラが、オンラインになるプロセスを開始する、請求項 1 に記載の方法。

【請求項 7】

前記オンラインになるプロセスは、前記第 2 のメモリの自己リフレッシュモードを終了させた後で、前記冗長コントローラシステムからの前記第 1 のコントローラの除去が完了する前に、直ちに開始されるように構成される、

請求項 6 に記載の方法。

【請求項 8】

請求項 1 から 7 のいずれかに記載の方法を実行するように構成された冗長コントローラシステム。

30

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は一般に、冗長コントローラを使用した冗長コントローラシステムおよびデータ格納システムに関し、さらに詳細には、オンラインコントローラ除去システムを有する冗長コントローラデータ格納システムおよび方法に関する。

【0002】

【従来の技術】

極めて信頼性の高い冗長データ格納システムを提供するために、多重コントローラシステムが使用されている。例えば、ハードディスクドライブ業界においては、多重コントローラシステムは、ディスクドライブの耐故障性およびディスクドライブ性能を向上させるために複数のディスクドライブを組み合わせ使用して使用する、RAID (redundant array of independent disks の略) システムの一部として使用されている。動作時には、RAID システムは、多重コントローラを使用することによって冗長性を持たせている。多重コントローラは、ユーザのデータを複数のハードディスク間でストライプしている。上記冗長アレイは、任意のコントローラから動作させることができる。多重コントローラが存在する場合、それらのコントローラを使用することにより性能が向上し、および / またはホストコンピュータシステムの接続ポート数が増加する。データにアクセスする場合、多重コントローラ RAID システムにより、すべてのハードディスクが同時に作動し、速度および信頼性が大幅に向上する。

40

50

【0003】

RAIDシステムの構成は、様々なRAIDレベルによって定義される。その様々なRAIDレベルの範囲は、データのストライピング（各ファイルのデータブロックを複数のハードディスクに分散して格納する）を提供し、冗長性を持たせることなくディスクドライブの速度および性能を向上させるLEVEL0から始まっている。RAID LEVEL1は、ディスクのミラリングを提供し、ミラーハードディスク対を介して、データの100%の冗長度をもたらしている（すなわち、データの同一ブロックが2つのハードディスクに書き込まれる）。他のディスクドライブにおけるRAID LEVELでは、様々なデータストライピングおよびディスクミラリングを提供し、性能を向上させるためのエラー修正、耐故障性、効率および/またはコストを改善している。

10

【0004】

RAID LEVEL5は、データをブロックに区切り、それらをディスクドライブに渡ってストライプしている。データブロックからパリティブロックが計算され、ディスクに格納される。すべてのデータブロックおよびパリティブロックは、異なるディスク上に格納される（ストライプされる）。ディスクドライブのいずれの1つが故障しても、データブロックまたはパリティブロックの1つが失われるだけである。その場合、冗長アレイは、失われたブロックを数学的に再生成することができる。また、RAID5は、データブロックおよびパリティブロックが格納されるディスクを循環させる（つまり、すべてのディスクは、そのディスク上にいくつかのパリティブロックを格納する）。RAID LEVEL6は、上記ステップをさらに進め、異なる数学式を用いて2つの「パリティ」ブロックを計算している。これにより上記冗長アレイは、2つの故障ディスクドライブを許容することができ、すべてのデータを再生成することができる。

20

【0005】

知られている多重コントローラシステムは、ミラー二重コントローラデータ格納システムを備えている。各コントローラはそれぞれ独自のメモリを含んでおり、そのほとんどは、「ミラーイメージ」すなわち他のメモリと同じ「メモリエメージ」である。ミラーメモリを二重コントローラで使用するにより、コントローラの1つまたはそのメモリが破損または失われた場合、速やかに復旧することができ、かつ、データの損失を防止することができる。メモリのミラーコピーがなければ、コントローラが突然故障した場合に、そのコントローラ上の重要なデータが失われてしまうであろう。

30

【0006】

例えば、コントローラAおよびコントローラBを有するミラーメモリ二重コントローラシステムでは、ミラー読出しおよび書込みにより、コントローラAのメモリが、コントローラBのメモリの「ミラーイメージ」になる。損失すなわちコントローラBが故障すると、すべてのシステムの動作は自動的にコントローラAに切り換わり、コントローラAがシステム全体をラン、すなわち動作させる。

【0007】

コンピュータシステムのアプリケーション数の増加により、プロセッサ故障時間の極めて大きな制限を始めとする、非常に高度の信頼性が要求されている。例えば、知られているシステムでは、コントローラの総故障時間は、1年当たり5分未満であることが要求されている。通常、損失すなわちコントローラの1つが故障すると、関連するデータ格納システムの冗長性および信頼性を維持するために、直ちにコントローラを交換しなければならない。上記の要求により、通常、高度の信頼性および「動作可能時間」が要求されるシステムの場合、交換するコントローラを、オンライン挿入すなわち「ホット」挿入しなければならない。その間、他のコントローラ（例えば、コントローラA）は、動作状態を維持していなければならない。オペレーティングシステムは、交換コントローラが挿入されたことを自動的に認識する。

40

【0008】

通常、多重コントローラシステムはホストに接続される。従ってホストシステムは、しばしばコントローラボードを交換することでかなりの長時間の間データ格納システムを停止

50

させたり、ホストシステムのタイムアウトを引き起こしたりしないよう要求する。

【 0 0 0 9 】

動作中のシステムに交換コントローラを挿入すると、交換コントローラが、動作中のシステムにおいて試験され加えられるまでの間、システム利用上の損失をしばしば引き起こす。ミラーメモリシステムの一部として交換コントローラが追加される場合、動作中のシステムに交換コントローラを追加することに関連して問題が増加する。

【 0 0 1 0 】

知られているミラーメモリ二重コントローラシステムでは、コントローラ A をシステム内で動作中に交換コントローラ B をホット挿入する場合、コントローラ A および交換コントローラ B の両方をリセットし、コントローラの各々に、プロセッサのサブシステムを自己診断させなければならない。

10

【 0 0 1 1 】

各コントローラは、自己の共用メモリシステムを試験し、ハードウェアが正常に動作していることを確認する。各コントローラは、その共用メモリの内容をチェックし、システムに対するメモリイメージの「有効」性を確認する。この例では、コントローラ A のみが、システムの有効メモリイメージを持つことになる。

【 0 0 1 2 】

次に、各コントローラは、互いのレビジョン、システムの最新ビュー、およびシステムが最後にアクティブであったときのシステム状態、に関する情報を交換する。これらの情報を共有した後、ファームウェアは、いずれのコントローラが有効メモリイメージを有しているかを決定する。この例の場合、コントローラ A が有効メモリイメージを有している。コントローラ A の共用メモリイメージは、コントローラ B にコピーされ検証されるが、そのためには、コントローラ A のプロセッサは、コントローラ A のすべての共用メモリを読み出し、コントローラ B のすべての共用メモリロケーションに書き込まなければならない。

20

【 0 0 1 3 】

次に、コピーオペレーションが成功したことを検証するために、両コントローラのメモリが読み出され、比較される。メモリシステムが大きい場合、上記のプロセスには数分の時間が必要である。最終構成ステップが実行されると、コントローラはオンラインになり、全動作状態になる。

30

【 0 0 1 4 】

上記プロセスの多くのステップは、実行するにあたり数十秒の時間が必要なことがある。コントローラ A の共用メモリイメージをコントローラ B にコピーし、検証するプロセスには、数分の時間が必要なことがある。ホット挿入に必要なこの延長期間の間に、ほとんどのホストコンピュータのオペレーティングシステムがタイムアウトすることになる。

【 0 0 1 5 】

【発明が解決しようとする課題】

システムの故障時間を短縮し、かつ、ホストコンピュータのオペレーティングシステムタイムアウトの原因になることのない冗長ミラーメモリ多重コントローラシステムに使用するためのホット挿入および/またはシステムおよび方法を持つことが望ましい。また、システムの故障時間またはホストのタイムアウトを最小にする、有効なコントローラリセット処理方法を持つことが望ましい。

40

【 0 0 1 6 】

【課題を解決するための手段】

本発明は、冗長コントローラを使用した多重コントローラシステムおよびデータ格納システムに関し、さらに詳細には、オンラインコントローラ除去システムを有する冗長コントローラデータ格納システムおよび方法に関する。

【 0 0 1 7 】

本発明の一実施形態によれば、冗長コントローラシステムからコントローラをオンライン除去する方法が提供される。上記冗長コントローラシステムは、第 1 のコントローラおよ

50

び第2のコントローラを備えており、上記方法には、冗長コントローラシステムからの、上記第1のコントローラの部分除去を検出するステップが含まれている。未処理のメモリアクセスの完了を含むシャットダウンシーケンスが、第1のコントローラおよび第2のコントローラに対して実行される。第1のコントローラは、自己リフレッシュモードに置かれる第1のメモリを有するように定義されている。冗長コントローラシステムからの第1のコントローラの除去は、上記第1のメモリによってモード化される自己リフレッシュが完了した後、終了する。

【0018】

本発明の他の実施形態によれば、冗長コントローラシステムからコントローラをオンライン除去する方法が提供される。上記冗長コントローラシステムは、第1のコントローラおよび第2のコントローラを備えている。第1のコントローラは、第1のプロセッサを備えるように定義され、第2のコントローラは、第2のプロセッサを備えるように定義されている。冗長コントローラシステムからの上記第1のコントローラの部分除去が検出される。上記第1のコントローラおよび第2のコントローラの未処理のメモリアクセスを完了させるための上記第1のプロセッサへの割込みおよび第2のプロセッサへの割込みを含むシャットダウンシーケンスが、第1のコントローラおよび第2のコントローラに対して実行される。第1のコントローラは、第1のメモリを有するように定義され、第1のメモリを自己リフレッシュモードにしている。冗長コントローラシステムからの第1のコントローラの除去は、上記第1のメモリによる自己リフレッシュモードが完了した後、終了する。

【0019】

本発明の他の実施形態によれば、冗長コントローラをオンライン除去するために構成された冗長コントローラシステムが提供される。システムは、第1のメモリを含む第1のコントローラ、第1のプロセッサ、および上記第1のコントローラの部分除去を早期検出するためのシステムを備えており、第1のコントローラの部分除去が検出されると、第1のコントローラは、上記第1のメモリに対して未処理になっているメモリアクセスの完了を含むシャットダウンシーケンスを実行する。第2のコントローラは、第2のメモリおよび第2のプロセッサを備えており、上記第1のコントローラの部分除去が検出されると、上記第1のメモリに対して未処理になっているメモリアクセスの完了を含むシャットダウンモードを実行する。第1のコントローラおよび第2のコントローラによるシャットダウンシーケンスが完了すると、上記第2のメモリが自己リフレッシュモードに置かれ、冗長コントローラシステムからの第1のコントローラの除去が完了する。

【0020】

【発明の実施の形態】

本発明の好ましい実施形態について、本発明を實踐することができる特定の実施形態の説明用として示す、本明細書の一部を形成する添付の図面に照らして、以下に詳細に説明する。本発明の範囲を逸脱することなく他の実施形態を利用し、また、構造あるいはロジックを変更することができることを理解しなければならない。したがって以下の詳細説明を、本発明を制限するものとして捕えてはならない。本発明の範囲については、特許請求の範囲の各クレームによって定義されている。

【0021】

図1は、本発明による冗長コントローラデータ格納システムの一例示的实施形態を、符号30で総括して示したものである。冗長コントローラデータ格納システムは、冗長ミラーメモリ、オンラインすなわち「ホット」挿入システムを有する多重コントローラシステム、及びシステム故障時間を短縮しかつコントローラの交換中にホストコンピュータのオペレーティングシステムがタイムアウトしないための方法を提供している。本発明の一態様では、冗長コントローラデータ格納システム30は、二重コントローラシステムである。本明細書において説明する例示的实施形態は、二重コントローラシステムを採用しているが、これらの実施形態は、他の多重コントローラ環境（つまり、3つ以上のコントローラを有するシステム）にも適用することができる。

【0022】

本発明によるコンポーネントは、マイクロプロセッサ、プログラマブルロジックすなわち状態マシン、およびファームウェアを介してハードウェアで実施することができ、あるいは所定の装置内のソフトウェアで実施することもできる。好ましい実施形態では、ソフトウェア中に本発明による１つまたは複数のコンポーネントが存在し、ハードウェアを介して使用されている。また、本発明によるコンポーネントは、１つまたは複数のコンピュータ可読媒体のソフトウェア中に存在させることもできる。本明細書で使用されているコンピュータ可読媒体という用語は、フロッピーディスク、ハードディスク、CD-ROM、フラッシュメモリ、読出し専用メモリ（ROM）、およびランダムアクセスメモリ（RAM）など、任意の種類の揮発性メモリまたは不揮発性メモリを含むものとして定義されている。また、本発明によるシステムは、注文製作装置ハードウェアおよび／または専用単一目的ハードウェアが組み込まれた、マイクロプロセッサ埋込みシステム／装置を使用することができる。

10

【0023】

一例示的实施形態では、システム30は、第1のコントローラ32および第2のコントローラ34を有する冗長ミラーコントローラデータ格納システムである。第1のコントローラ32および第2のコントローラ34は、（例えば）ディスクアレイなど、データ格納システム36へのデータの冗長又はミラーによる読出し及び書込み用に構成されているが、このデータ格納システム36は通信バス38を介しており、（例えば）RAID LEVEL 1においてはここで、ミラー書込みや、アレイがコピーからの読み出しを行うモードを含み、RAID LEVEL 5または6の場合、ユーザのアクセスは、上記ディスクアレイ間でストライプされる。さらに、第1のコントローラ32および第2のコントローラ34は、通信バス40を介して互いに通信している。第1のコントローラ32および第2のコントローラ34は、通信バスプロトコルを使用してデータ格納システム36と通信し、かつ、相互に通信している。一態様では、上記通信バスプロトコルは標準プロトコルである。他の適切な通信バスプロトコルは、本出願から当分野の技術者には明らかになるであろう。データ格納システム36は、磁気ハードディスクデータ格納システムを備えている。他の態様では、データ格納システム36は、フラッシュメモリ、ランダムアクセスメモリ（RAM）、CD書込み可能媒体、光磁気媒体等、他の読出し／書込み可能データ格納媒体を備えている。

20

【0024】

冗長コントローラデータ格納システムは、ホストまたはコントロールシステムインタフェース42を介して、ホストまたはコントロールシステムと通信するように構成されている。ホストまたはコントロールシステムは、サーバであり、コンピュータネットワークであり、中央計算機であり、あるいは他のコントロールシステムである。一態様では、冗長コントローラデータ格納システム30は、ホストとインタフェースし、RAIDシステム（例えばRAID LEVEL 0、RAID LEVEL 1、RAID LEVEL 2、RAID LEVEL 3、RAID LEVEL 4、RAID LEVEL 5、またはRAID LEVEL 6システム）として動作するように構成される。

30

【0025】

一実施形態では、第1のコントローラ32は、「ミラーされた」メモリ50、タスクプロセッサ52、およびシステムオペレーションプロセッサ54を備えている。同様に、第2のコントローラ34も、「ミラーされた」メモリ56、タスクプロセッサ58、およびシステムオペレーションプロセッサ60を備えている。第1のコントローラ32および第2のコントローラ34は、メモリ50およびメモリ56をミラーメモリシステムの一部として動作させる「メモリコントローラ」を備えている。本明細書で使用する「ミラーメモリ」という用語は、１つのメモリのメモリエイジが他のメモリに複製すなわち「ミラーされている」システムを含むものとして定義されている。本発明においては、第1のコントローラ32のメモリ50は、第2のコントローラ34のメモリ56に複製すなわち「ミラーされている」。二重コントローラミラーメモリシステムは、冗長コントローラシステム30に故障許容環境を提供しており、コントローラ的一方、あるいはコントローラメモ

40

50

リシステムの一方が故障した場合に、もう一方のコントローラおよびそのミラーメモリが存在することにより、中断することなく故障が復旧され、システムコマンドの処理が継続される。また、コントローラの一方を除去し、かつ、挿入する場合、本発明により、もう一方のコントローラを介してオペレーティングシステムが維持され、ホストがタイムアウトする期間以内にシステム故障時間が短縮される。ミラーメモリ二重コントローラディスク格納システムの一例示的实施形態が、本出願人に譲渡された、1997年12月16日発行の米国特許第5,699,510号に開示されている。また、同じくこの出願の出願人に譲渡された、1999年7月27日発行の米国特許出願第5,928,367号に、他のミラーメモリ二重コントローラディスク格納システムが開示されている。

【0026】

本発明による冗長コントローラデータ格納システム30では、各コントローラ32および34は、それぞれ独自のメモリ50および56を備えており、メモリは「ミラーイメージ」、すなわち上で示したように、他のメモリと同じ「ミラーイメージ」を有している。ミラーメモリは、一方のコントローラまたはそのメモリが故障すなわち損失した場合における迅速な復旧を可能にしている。一態様では、ミラー読み出しおよびミラー書き込みにより、第1のコントローラ32のメモリ50が、第2のコントローラ34のメモリ56の「ミラーイメージ」になる。第2のコントローラ34が損失すなわち故障すると、すべてのシステムオペレーションは、第1のコントローラ32に自動的に切り換えられ、システムに他のコントローラが挿入されるまでの間、単一コントローラシステムとして第1のコントローラ32がシステム全体をラン、すなわち動作させる。

【0027】

本発明による冗長コントローラデータ格納システム30により、一方のコントローラのホット挿入時、すなわちオンライン挿入時における冗長コントローラシステムの継続動作が提供される。例えば、第2のコントローラ34が損失すなわち故障すると、冗長コントローラシステムは、第1のコントローラ32を介して動作する。第2のコントローラ34は、システム30にオンライン挿入すなわちホット挿入することができる。詳細には、システムに第2のコントローラ34をもたすためのホット挿入プロセスの間、システムオペレーションプロセッサ54が、データ格納システム36に対する、メモリ50を介したデータ読み出しおよび書き込みなどのシステムオペレーションコマンドの処理を継続する。タスクプロセッサ52は、システムオペレーションプロセッサ54がシステムオペレーションコマンドを処理している間、冗長コントローラデータ格納システム30を遅延させることなく、バックグラウンドタスクを処理する。

【0028】

好ましい一実施形態では、タスクプロセッサ52は、システムオペレーションプロセッサ54がシステムオペレーションコマンドの処理を継続している間、第1のミラーメモリ50のメモリイメージを第2のメモリ56にコピーするように動作している。したがって、第2のコントローラ34を冗長コントローラシステム30にホット挿入することにより、システムオペレーションコマンドの処理が遅延することはない、および/または、ホストシステムインタフェース42を介したホストシステムがタイムアウトすることもない。一例示的实施形態では、タスクプロセッサ52は、システムオペレーションプロセッサ54または他のシステムプロセッサを直接煩わせることなく、専用データ処理ハードウェアを介してバックグラウンドタスクを実行している。一態様では、上記データ処理ハードウェアは、専用集積回路(ASIC)の一部として、インテリジェントDMAエンジンに結合されている。タスクプロセッサ52は、第1のコントローラ32がシステムオペレーションプロセッサ54を介して動作を継続している間、特定のバックグラウンドタスクを処理する能力を有している。一態様では、タスクプロセッサ52は、メモリ間コピータスク、メモリ自己診断、およびその他のタスクを実行している。

【0029】

図2は、一括して80で示される、本発明による冗長コントローラデータ格納システムにコントローラをホット挿入する方法の一例示的实施形態を示したものである。上記方法に

10

20

30

40

50

は、第 1 のメモリ、タスクプロセッサ、およびシステムオペレーションプロセッサを備えるように第 1 のコントローラを構成するステップが含まれている。上記第 1 のメモリは、第 1 のメモリイメージを含んでいる。図に示す一例示的实施形態では、第 1 のコントローラ 32 は、メモリ 50、タスクプロセッサ 52、およびシステムオペレーションプロセッサ 54 を備えるように構成されている。冗長コントローラシステム 30 は、84 に示すように、単一コントローラシステムとして、第 1 のコントローラ 32 を介して動作する。85 で、システムオペレーションコマンドが、システムオペレーションプロセッサ 54 を介して処理される。86 で、冗長コントローラシステム 30 に第 2 のコントローラ 34 が挿入される。第 2 のコントローラ 34 は、第 2 のメモリ 56 を備えている。88 で、第 1 のコントローラを介してシステムオペレーションコマンドが処理されている間に、タスクプロセッサ 52 を用いてバックグラウンドタスクが処理される。上記バックグラウンドタスクには、メモリ 50 の第 1 のイメージの、第 2 のメモリ 56 へのコピーが含まれている。

【0030】

図 3 は、一括して 90 で示される、本発明による冗長コントローラデータ格納システムにコントローラをホット挿入する方法の他の例示的实施形態を示したものである。上記方法には、一括して 92 で示される、第 1 のメモリ 50、タスクプロセッサ 52、およびシステムオペレーションプロセッサ 54 を備えるように第 1 のコントローラ 32 を構成するステップが含まれている。第 1 のメモリ 50 は、第 1 のメモリイメージを含んでいる。94 で、冗長コントローラシステム 30 は、第 1 のコントローラ 32 を介して動作する。96 で、システムオペレーションコマンドが、システムオペレーションプロセッサ 54 を介して処理される。98 で、冗長コントローラシステム 30 に第 2 のコントローラ 34 が挿入される。第 2 のコントローラ 34 は、第 2 のメモリ 56 を備えている。100 で、第 1 のメモリ 50 は、第 2 のメモリ 56 へのミラー書込み用、及び共用すなわちミラーメモリ 50 に対するローカル読取り専用として構成される。従って第 1 のコントローラ 32 は、その独自のメモリイメージを読み出すために動作することができるが、第 2 のコントローラ 34 が、冗長コントローラデータ格納システム 30 内で完全に動作状態になるまで（すなわち自己診断が終了し、オンラインになるまで）、ミラー書込みおよびミラー読出しとして動作することはない。102 で、第 1 のコントローラ 32 を介してシステムオペレーションコマンドが処理されている間に、バックグラウンドタスクが処理される。バックグラウンドタスクは、タスクプロセッサ 52 を使用して処理される。バックグラウンドタスクには、メモリ 50 の第 1 のイメージの、第 2 のメモリ 56 へのコピーが含まれている。

【0031】

図 4 は、本発明による冗長コントローラデータ格納システムの他の例示的实施形態を、一括して 110 で示したものである。冗長コントローラデータ格納システム 110 は、本明細書において既に説明した冗長コントローラデータ格納システム 30 と類似している。冗長コントローラデータ格納システム 110 は、ホストシステムのタイムアウトの原因になり得る、システムオペレーティングコマンドの処理に対するあらゆる割込みを最小化する、冗長コントローラデータ格納システムへのコントローラのホット挿入システムおよび方法を備えている。

【0032】

冗長コントローラデータ格納システム 110 は、第 1 の冗長コントローラ 112 および第 2 の冗長コントローラ 114 を備えている。第 1 のコントローラ 112 は、第 1 のミラーメモリ 120、第 1 のメモリコントローラ 122、および第 1 のシステムオペレーションプロセッサ 124 を備えている。一態様では、第 1 のコントローラ 112 は、ディスクインタフェース 126 およびディスクインタフェース 128 を介してデータ格納システムと通信し、また、ホストインタフェース 130 を介して、ホストまたはコントロールシステムと通信している。一態様では、第 1 のコントローラ 112 は、通信バス 132 を介してディスクインタフェース 126、ディスクインタフェース 128、およびホストインタフェース 130 と通信している。一実施形態では、通信バス 132 は、当分野の技術者に知られている PCI バスとして構成されている。一実施形態では、130、126 および 1

10

20

30

40

50

28で示すホストインタフェースおよびディスクインタフェースは、「FCループ」として動作することができるファイバチャネルバスである。他の適切なバス構成については、本出願を読めば、当分野の技術者には明らかになるであろう。

【0033】

一態様では、メモリコントローラ122は、タスクプロセッサ134、割込みロジック136、およびメモリバッファ/通信モジュール138を備えている。一態様では、タスクプロセッサ134は、専用ファームウェアおよび/またはメモリバッファコンポーネントを備えており、システムオペレーションプロセッサ136を介したシステムオペレーションコマンドの処理を中断することなく、定義済みバックグラウンドタスクを処理している。メモリコントローラ122のためのホットプラグ警告/早期検出システムが、142で示されている。同様に、メモリコントローラ122のためのリセットロジックが、140で示されている。

10

【0034】

同様に、第2の冗長コントローラ114は、第2の共用すなわちミラーメモリ160、第2のメモリコントローラ162、および第2のシステムオペレーションプロセッサ164を備えている。第2のコントローラ114は、ディスクインタフェース166およびディスクインタフェース168を介してデータ格納システムと通信し、また、ホストインタフェース170を介して、ホスト/コントロールシステムと通信している。第2のコントローラ114は、通信バス172を介してディスクインタフェース166、ディスクインタフェース168、およびホストインタフェース170と通信している。

20

【0035】

第2のメモリコントローラ162は、タスクプロセッサ174、割込みロジック176、およびメモリ/通信モジュール178を備えている。第2のコントローラ114のためのリセットロジックは、180で示されている。ホットプラグ警告/早期検出システムが第2のメモリコントローラ162に設けられており、182で示されている。第1のコントローラ112および第2のコントローラ114は、両コントローラ間の通信バスを介して通信している。一態様では、ミラーバス200が、第1のコントローラ112の第1のメモリコントローラ122と、第2のコントローラ114の第2のメモリコントローラ162とをリンクしている。また、交互通信経路が、第1のメモリコントローラ122と第2のメモリコントローラ162の間に設けられており、202で示されている。交互通信経路202は、第1のメモリコントローラ122のメモリ/通信モジュール138にリンクされ、かつ、第2のメモリコントローラ162の第2のメモリ/通信モジュール178にリンクされている。存在検出ライン204は、コントローラの存在について第1のコントローラ112と第2のコントローラ114の間の通信を提供している（例えば、ホット挿入プロセスの一部として）。

30

【0036】

一実施形態では、メモリ120およびメモリ160は、ランダムアクセスメモリ(RAM)である。一例示的实施形態では、上記ランダムアクセスメモリは、同期ダイナミックRAM(SDRAM)である。一態様では、メモリ120およびメモリ160のサイズは、512バイトから数ギガバイトの範囲に及んでいる。好ましい一実施形態では、メモリ120およびメモリ160は、バッテリーバックアップ式RAMなどの不揮発性メモリであり、そのため、電源がシャットダウン（例えばコントローラリセット）された場合においても、メモリはその記憶内容（すなわち、記憶状態）を保持している。

40

【0037】

一態様では、メモリコントローラ122およびメモリコントローラ162は、特定用途向け専用集積回路(ASIC)チップまたはモジュールの一部である。タスクプロセッサ134およびタスクプロセッサ174は、システムオペレーションプロセッサを介してシステムコマンドが処理されている間に、定義済み専用バックグラウンドタスクを処理するように動作する。

【0038】

50

一実施形態では、タスクプロセッサ134によって実行されるすべてのバックグラウンドタスクすなわち機能は、メモリ120およびメモリ160に格納されているデータに対して働き、これらのタスクの結果は、適当なメモリ120またはメモリ160に置き戻される。タスクプロセッサ134およびタスクプロセッサ174は、専用データ処理ハードウェアを使用して、プロセッサ124あるいはプロセッサ164などの他のシステムプロセッサを直接煩わせることなく、バックグラウンドタスク機能を実行する。一態様では、タスクプロセッサ134および/またはタスクプロセッサ174は、インテリジェントDMAエンジンに結合された、ASICチップまたはモジュールの一部であるデータ処理ハードウェアを利用している。タスクプロセッサ134およびタスクプロセッサ174の例示的实施形態のさらに詳細については、本明細書の中で後述する。

10

【0039】

一態様では、第1のメモリ120と第2のメモリ160の間のミラー読出しおよびミラー書込みは、ミラーバス200を介して達成される。さらに、交互通信経路すなわちバス202が、メモリコントローラ122とメモリコントローラ162の間に設けられている。したがって、冗長コントローラデータ格納システムにコントローラが挿入され、冗長コントローラシステムの一部として未だオンラインになっていない状態において、第1のメモリコントローラ122は、交互通信経路202を介して第2のメモリコントローラ162と通信することができる。このような通信には、ハードウェアおよびファームウェアのレビジョン情報の交換、システム内でコントローラが交換された場合に検出するシリアル番号の交換、互いの動作状態に関する情報の交換、およびホット挿入シーケンスにおける次のステップへの移行のタイミングの相互通知が含まれている。また、故障によりミラーバス200を介した通信が妨害された場合に、通信バスを用いて、いずれのコントローラが動作状態を維持すべきかが交渉される。ファームウェアの他の領域は、他の目的のためにこの通信バスを使用している。ホットプラグ警告142およびホットプラグ警告182は、コントローラがホット挿入されていることを冗長コントローラシステムに知らせる早期検出信号を、対応するメモリコントローラ122およびメモリコントローラ162に提供するように動作する。早期警告ロジックは、リセットロジックと共同して、ホット挿入されるコントローラがオンラインになるまで、コントローラをリセット状態に保持する。コントローラがホット除去されている間、早期検出信号は、コントローラが除去されていることについての早期警告を提供する。ホットプラグ警告142およびホットプラグ警告182早期検出システムは、機械的手段あるいは電気的手段、例えばコネクタピン、押しボタン警告、センサ検出（例えば光センサ）、またはその他の検出システムの使用を介して、早期検出信号を受け取ることができる。存在検出ライン204は、コントローラがシステムから除去されたこと、あるいはシステムに挿入されたことを他のコントローラに知らせるように動作する。

20

30

【0040】

プロセッサ124およびプロセッサ164は、メモリコントローラ122およびメモリコントローラ164を介して、対応するメモリ120およびメモリ160と通信する、システムコマンドを操作するためのシステムオペレーションプロセッサである。このようなシステムコマンドには、ディスクインタフェース126、132、166および168を介して、対応するデータ格納システムからデータを読出し、かつ、システムにデータを書込むための、ホストインタフェース130およびホストインタフェース170を介して受け取るシステムコマンドが含まれている。プロセッサ124およびプロセッサ164は、システム割込みオペレーション、リセットオペレーション、または他のシステムプロセスの処理および管理など、他のシステムオペレーションを実行するように動作する。

40

【0041】

図5は、本発明による冗長コントローラシステムに使用されるタスクプロセッサの一例示的实施形態を示したものである。タスクプロセッサ134は、一例として示したものであるが、タスクプロセッサ174についても、タスクプロセッサ134と同様である。タスクプロセッサ134は、データ処理ハードウェアを介して、定義済み機能を実行すること

50

が好ましい。これらのタスクは、「バックグラウンドタスク」として処理されるため、システムオペレーションプロセッサ 124 を介してシステムコマンドが動作している間に完了させることができる。一例示的实施形態では、タスクプロセッサ 134 は、メモリイメージをメモリ 120 とメモリ 160 の間でコピーするためのメモリ間コピータスク 206 を含んでいる。また、タスクプロセッサ 134 は、関連するメモリ 120 の自己診断を実行するための、1 つまたは複数のメモリ自己診断タスク 208 を含んでいる。メモリ自己診断 208 は、冗長コントローラシステムへのコントローラの挿入時、あるいは冗長コントローラシステム 110 が動作中の任意のタイミングで実行させることができる。典型的なメモリ自己診断には、メモリイメージ、メモリチャンクすなわちデータブロックの読み出し、および内部バッファ（例えば、メモリコントローラ 122 の内部バッファ）への保管が含まれている。テストパターンがメモリブロックに書き込まれ、それがリードバックされて正当性が検証される。このステップは、多くのテストパターンを用いて繰り返される。一態様では、タスクプロセッサは、1 回のテストで 1 ないし 30 のパターンを実行することができる。内部バッファに格納されていた元のデータブロックが、外部メモリブロックにライトバックされる。このプロセスは、メモリのすべてのブロックのテストが終了するまで繰り返される。タスクプロセッサ 134 の他のタスクには、二重ブロックパリティ生成 210、単一ブロックパリティ生成 212、ブロックパターン認識 214、およびチェックサム生成 216 が含まれている。

【0042】

図 6 は、タスクオペレーションを処理するために、タスクプロセッサ 134 およびタスクプロセッサ 174 によって使用されるデータ構造の一例示的实施形態を示したものである。他の適切なデータ構造については、本出願を読めば、当分野の技術者には明らかになるであろう。一例示的实施形態では、要求元プロセッサが、タスク記述ブロック (TDB) をメモリ 120 に書き込んでいる。タスク記述ブロックには、コマンドコードおよび要求を処理するために必要なコマンド専用情報（ブロックアドレス、ブロックサイズ、データパターン、パリティ係数に対するポインタ等）が含まれている。要求元プロセッサは、次に、タスクプロセッサ 134 に対して局部的に、要求エントリを要求キュー（例えば、220 で示されるキュー「0」）に挿入する。このエントリには、コマンドコード要求ヘッダ 222、関連するタスク記述ブロックに対する TDB ポインタ 224、およびキューポインタ 226 で示される、応答のためのキュー番号が含まれている。

【0043】

キュー「0」220 信号が、空ではない信号を送ると、タスクプロセッサは、キュー 220 から要求エントリを読み出す。タスクプロセッサは、要求情報を用いてタスクデータブロック 228 を読み出し、読み出したタスクデータブロック 228 が要求と一致していることを確認する。次に、タスクプロセッサは、所望の機能を実行する。タスクプロセッサは、完了応答エントリを、230 で示される指定応答キューに置き、応答キュー 230 を通して要求元プロセッサ 124 に完了通知される。

【0044】

図 7 は、タスクプロセッサ 134 による処理に適したメモリブロックに分割された第 1 のメモリ 120 に含まれるメモリイメージの一例示的实施形態を示したものである。詳細には、タスクプロセッサ 134 によって処理されるバックグラウンドタスクは、特定のタスクプロセッサ 134 を備えるメモリコントローラ 122 の内部にバッファリングするには大き過ぎる、メモリ 120 に格納されているデータブロックに対して働くことができる。したがってタスクプロセッサ 134 は、メモリイメージすなわちブロックを、タスクプロセッサによって処理することができるサイズに対応するメモリブロックすなわちチャンク中に構成するように動作する。図に示す例示的实施形態では、第 1 のメモリ 120 に格納されているメモリイメージは、メモリブロック「1」232、メモリブロック「2」234、メモリブロック「3」236、メモリブロック「4」238 ないしメモリブロック「N」240 中に構成される。一態様では、各チャンクは最大 512 バイトであり、メモリコントローラ 122 の内部に内部バッファリングするには十分に小さく、かつ、タスク処

10

20

30

40

50

理システムを有効に使用するには十分な大きさである。一態様では、タスクプロセッサ 1 3 4 は、メモリコントローラ 1 2 2 の制限内で動作しつつ、メモリイメージ当たりのメモリブロック数が最小になるメモリブロックサイズを構成するように動作している。最大使用可能メモリブロックが 5 1 2 バイトである一態様では、タスクプロセッサ 1 3 4 は、ブロック中にメモリイメージを構成し、第 1 のメモリブロックおよび最後のメモリブロックのみ、最大すなわち 5 1 2 バイト未満にすることができる。図に示す例示的实施形態では、メモリブロック「1」2 3 2 およびメモリブロック「N」2 4 0 を、最大メモリブロックサイズ未満にすることができる。メモリブロック「2」2 3 4、メモリブロック「3」2 3 6、およびメモリブロック「4」2 3 8 等は、最大メモリブロックサイズ（例えば 5 1 2 バイト）にすることができる。

10

【0045】

タスクプロセッサ 1 3 4 およびシステムオペレーションプロセッサ 1 2 4 は、いずれもメモリ 1 2 0 に格納されているデータに対して動作する。冗長コントローラシステム 1 1 0 は、冗長コントローラシステムへの第 2 のコントローラの追加を含め、タスクプロセッサ 1 3 4 を介してタスクが処理されている間、システムコマンドの処理を継続することができるように構成されることが望ましい。したがって優先順位は、タスクプロセッサ 1 3 4 と他のシステムオペレーションとの間で割り振られるが、これにはメモリ 1 2 0 にアクセスするためのプロセッサ 1 2 4 を介して達成されるものなどが挙げられる。好ましい一実施形態では、タスクプロセッサ 1 3 4 に割り当てられている優先順位は、オペレーティングシステムの性能が、タスクプロセッサ 1 3 4 を介したバックグラウンドタスクの動作によって極端に低下しないように、プロセッサ 1 2 4 よりも低い優先順位（例えば、最下位の優先順位）になっている。別法としては、タスクプロセッサ 1 3 4 のメモリアccess優先順位を、他のシステムオペレーションと同じか、あるいはそれより高くすることができる。あるいは、ファームウェアを利用して、タスクプロセッサ 1 3 4 を介して達成される個々のタスクのメモリアccess優先順位を高くすることができる。

20

【0046】

図 8 ないし図 1 0 は、システムの中断を最小にする、本発明による冗長コントローラシステムへのコントローラのオンライン「ホット」挿入の一例示的实施形態を示したもので、ここでも、本明細書において既に説明した図 1 ないし図 7 が参照されている。

【0047】

図 8 に、本発明による冗長コントローラシステムにコントローラをホット挿入する方法の一例示的实施形態を示す図が、一括して 2 5 0 で示されている。この例示的实施形態では、冗長コントローラシステムは、第 1 のコントローラ 1 1 2 を介して動作しており、第 2 のコントローラは、冗長コントローラシステムから除去されている。2 5 2 で、冗長コントローラシステム 1 1 0 は、単一コントローラシステムとして、第 1 のコントローラ 1 1 2 を介して動作する。2 5 4 で、第 1 のコントローラ 1 1 2 が、冗長コントローラシステム 1 1 0 に第 2 のコントローラ 1 1 4 が追加されたことを検出する。冗長コントローラシステム 1 1 0 に第 2 のコントローラ 1 1 4 が追加されたことを検出した後、第 1 のコントローラ 1 1 2 は、単一コントローラシステムとしての動作を継続する。第 1 のコントローラは、冗長コントローラシステム 1 1 0 に第 2 のコントローラ 1 1 4 が追加されたことを示す検出信号を受け取る。一態様では、冗長コントローラシステム 1 1 0 に第 2 のコントローラ 1 1 4 がホット挿入される場合、ホット挿入されるコントローラは、コントローラが完全に挿入され、所定の位置にラッチされるまで、リセット状態に保持される。存在検出ラインが、新しいコントローラの装着を検出し、挿入されたコントローラ 1 1 4 の存在が、存在検出ライン 2 0 4 を介して第 1 のコントローラ 1 1 2 に通信される。

30

40

【0048】

2 5 6 で、第 2 のコントローラ 1 1 4 に電源が投入され、所定の位置へのラッチ待ち状態になり、次に、プロセッササブシステムの自己診断が実行される。プロセッササブシステムの自己診断には、フラッシュ ROM に存在しているプロセッササブシステムのファームウェアイメージのテスト、およびマイクロプロセッサローカルメモリのテストが含まれて

50

おり、周辺チップレジスタおよびデータ経路のテストが実行される。258で、第1のコントローラ112および第2のコントローラ114が、交互通信経路202を介して通信する。第2のコントローラ114が冗長コントローラシステム110の一部として未だ「オンライン」になっていない状態であっても、第1のコントローラ112および第2のコントローラ114は、交互通信経路202、メモリ/通信モジュール138、およびメモリ/通信モジュール178を介して、相互に通信する。交互通信経路202を介した、第1のコントローラ112と第2のコントローラ114の間のサンプル通信には、コントローラ間の互換性を確認するためのハードウェアおよびファームウェアのレビジョン情報の交換、テストの結果が出てテストが完了した時の通知、及びホット挿入プロセス中におけるコントローラ間の同期点の通信が含まれている。

10

【0049】

260で、第2のコントローラ114が、その共用メモリに対する自己診断の実行を含む自己診断を継続する。本明細書において既に説明したように、これらのテストを、バックグラウンドタスクとしてタスクプロセッサ174を介して、システムを中断させることなく実行することができる。262で、すべてのテストに合格しない場合、リカバリモード264に入る。リカバリモード264には、冗長コントローラシステムおよび/またはホストへのエラー状態の提供が含まれている。一実施形態では、コントローラに不良のマークが施され、オフラインに維持される。もう一方のコントローラに対して、上記プロセスが繰り返される。すべてのテストに合格した場合は、266で、第2のコントローラ114が、第1のコントローラ112にメッセージを送り、冗長コントローラシステム110に追加されるべく準備が整ったことを知らせる。

20

【0050】

図9に、本発明による冗長コントローラシステムにコントローラをホット挿入する方法をさらに示す図が、一括して220で示されている。272で、第1のコントローラ112は、メモリ120を共用書込みおよびローカル読み出し専用として位置づける。したがって、その時点以降、メモリ120に書き込まれるデータはすべて、同時に第2のコントローラ114のメモリ160にミラーすなわち書き込まれる。この時点では、メモリ160のメモリエイメージは、メモリ120のメモリエイメージの「ミラー」コピーではないため、冗長コントローラシステム110の一部として、メモリ120から読み出すことができるのはデータのみである。274で、第1のコントローラ112の許可が出るまで（例えば、交互通信経路202を介して）、第2のコントローラ114による共用メモリ120および独自のメモリ160への書込みが禁止される。

30

【0051】

276で、第1のコントローラ112が、バックグラウンドタスクを介して、その共用メモリ120のすべてを、共用メモリ120の同一ロケーションにコピーバックする。詳細にはタスクプロセッサ134はバックグラウンドタスクを含んでいるが、ここでタスクプロセッサ134は、メモリ120からメモリブロックを読み出し、読み出したメモリブロックをバッファに格納し、そしてメモリ120の同一ロケーションに書き込んでいる。第1のコントローラ112は、共用書込みモードに構成されているため、上記オペレーションの結果、第1のコントローラ112は、メモリブロック120からメモリブロックをローカルに読み出すことになるが、第1のコントローラ112が、メモリ120の同一ロケーションにライトバックする場合、同時に第2のコントローラ114のメモリ160の同一ロケーションに書き込むことになる。278で、上記バックグラウンドタスクの間、第1のコントローラ112は、プロセッサ124を介してシステムオペレーションコマンドを実行する動作状態を継続する。一態様では、メモリエイメージは、一度に1つのメモリブロックがコピーされる。バックグラウンドタスク完了後、ここで第1のメモリ120のメモリエイメージは第2のメモリ160のメモリエイメージのミラーであり、冗長コントローラシステム110への第2のコントローラ114の追加プロセスが、280で示す通り継続される。

40

【0052】

50

図10は、本発明による冗長コントローラシステムへのコントローラの追加方法の一例示の実施形態を示したものであり、290で示す通り、第1のコントローラ112のメモリイメージが第2のコントローラメモリ160のメモリにミラーされた後又はコピーされた後のものである。292で、第1のコントローラ112が再構成され、ミラー書込みおよびミラー読み出しモードに構成される。ここで第1のメモリ120および第2のメモリ160の両メモリに対して、データロケーションの読み出し及び書き込みが可能となる。294で、第1のコントローラ112がすべてのメモリロケーションを読み出し、タスクプロセッサ134を介したバックグラウンドタスクを用いて、第1のコントローラの共用メモリ120と第2のコントローラの共用メモリ160の一致性が比較される。したがって、この時点でシステムオペレーションコマンドが中断されることはない。

10

【0053】

296で、メモリが一致しない場合、リカバリモード298に入る。メモリが一致した場合は、300で、第1のコントローラ112および第2のコントローラ114が再構成され、第2のコントローラが冗長コントローラシステム110に追加される。この時点で冗長コントローラシステム110は、ミラーメモリ冗長コントローラシステムとして完全な動作状態になり、また、第2のコントローラは、ホストをタイムアウトさせることなく、システムオペレーションの処理を最短時間だけ中断させるだけで、冗長コントローラシステムにホット挿入される。

【0054】

コントローラのリセット

20

知られている二重コントローラシステムでは、一方のコントローラがリセットされると、もう一方のコントローラのマイクロプロセッサへの割込みが生じる。その場合、割り込まれた側のコントローラのプロセッサは、割込みおよびリセットの原因を処理しなければならない。この過去の方法には多くの欠点が知られている。第1のコントローラと第2のコントローラの間のリセットのタイミングは様々である。一方のコントローラが「リセットループ」に入ると、そのリセットは他のコントローラのプロセッサに割込みをかけるが、ここでコントローラシステムの性能に影響を及ぼす拡散動作を備え、洗練度の劣る、ミラリングインタフェースの状態変化に連動する。第2のコントローラが、「スタック」すなわちエラティックプロセッサを有している場合、第2のコントローラのプロセッサは、割込みを処理するためには役に立たないため、第1のコントローラにより、第2のコントローラをリセットさせることはできない。したがって第1のコントローラは、システムを復旧させるためには、リセットを発生する「ウォッチドッグ」モードに頼らざるを得ない。第2のコントローラが、システムを使用して格納されているデータを損傷させる原因となるためには、はるかに長い適期が存在している。

30

【0055】

図11および図12は、本発明による冗長コントローラシステムを使用してコントローラリセットを処理するシステムおよび方法の一例示の実施形態を示したものである。本発明による冗長コントローラシステムを使用したコントローラリセット処理方法により、ミラリングバスが使用可能状態になった場合に、リセットのみが第2のコントローラに伝播することができるように、一方のコントローラに対するリセットがローカライズされる。このローカライズにより、故障したコントローラが他のすべてのコントローラをリセット状態に保留することを回避している。図11および図12も、本明細書において既に説明した図1ないし図10が参照されている。

40

【0056】

図11に、本発明による多重冗長コントローラシステムにおけるコントローラリセット処理方法の一例示の実施形態が400で示されている。冗長コントローラシステムは、アクティブに接続されミラー対として動作する、第1のコントローラ112および第2のコントローラ114を備えており、402で、第2のコントローラ114に対するリセット状態が検出される。第2のコントローラ114はリセットされ、シャットダウンプロセスが開始される。404で、シャットダウンプロセスの一環として、第1のコントローラ11

50

2と第2のコントローラ114の間の通信リンクを介して、第1のコントローラ112に、コントローラがリセットされたことが通知される。好ましい実施形態では、ミラーバス200を介して、第1のコントローラ112に、第2のコントローラがリセットされていることが通知される。406で、第1のコントローラ112および第2のコントローラ114の両コントローラに対して、シャットダウンプロセスが実行される。したがって第1のコントローラ112および第2のコントローラ114は、同時にシャットダウンプロセスを通過する。

【0057】

図12に、本発明による二重コントローラシステムにおけるコントローラリセット方法の一例示的实施形態をさらに示す図が410で示されている。412で、両コントローラに対するシャットダウンプロセスが完了すると、第1のコントローラ112および第2のコントローラ114に電源が投入される。シャットダウンプロセスにより、すべての内部バッファがフラッシュされ、メモリが置かれる。414で、シャットダウンプロセスの一部として、第1のコントローラおよび第2のコントローラがリセットされる。416で、第1のコントローラ112と第2のコントローラ114の間のミラーバス200インタフェースの使用が禁止される。一態様では、リセットの作用により、第1のコントローラ112および第2のコントローラ114が、ミラーバス200の使用を禁止している。ミラーバスの使用を禁止する動作により、ミラーバスが再び使用可能になるまで、ボード間でのさらなるリセットの伝播を防止している。

【0058】

418で、第1のコントローラ112および第2のコントローラ114の各々が、自己診断を実行する。通常、自己診断には、内部メモリコントローラASIC122および162のテスト、およびすべてのデータ経路バスのテスト以外に、本出願において既に説明したように、マイクロプロセッササブシステムおよびSDRAMメモリのテストが含まれている。

【0059】

この時点では、第1のコントローラ112と第2のコントローラ114の間のミラーバス200インタフェースは、使用禁止の状態が維持されている。したがってコントローラの一方によって生じる、あるいは自己診断の結果として生じるあらゆるリセットまたは割込みによって、もう一方のコントローラが影響されることはない。420で、自己診断に合格しなかった場合、自己診断が不合格となったコントローラは、422のリカバリモードに入る。通常リカバリモードには、ホストコンピュータシステムに対するエラーの発生、故障の通知、及びアレイ中の「不良」コントローラが使用禁止されないよう除去することが含まれている。残った「正常」コントローラは、データのミラーが必要ない単一コントローラモードで動作を開始する。

【0060】

第1のコントローラ112および第2のコントローラ114の自己診断に合格すると、424で、第1のコントローラ112と第2のコントローラ114の間のミラーバス200インタフェースが使用可能になる。したがって第1のコントローラ112および第2のコントローラ114は、リセットまたはリセットの原因が除去されたことを確認し、ミラー対としての動作を、再び継続することができる。

【0061】

本発明による冗長コントローラシステムを使用した上述のコントローラリセット処理方法によると、ミラーバスと共に使用可能状態にされない限りリセットが第2のコントローラに伝播できないように、1つのコントローラに対するリセットをローカライズする。このローカライズにより、1つのコントローラに、システム内の他のすべてのコントローラをリセットさせることのないシステムファームウェアのための方法が可能となる一方で、多重コントローラボード間のリセットの同期化を、ハードウェアで管理する利点をもたらされる。

【0062】

コントローラの除去

冗長コントローラデータ格納システムからオンラインコントローラを除去する場合、コントローラが「部分除去」されるかあるいは正しく除去されない場合、故障時間が発生することになる。これは、通常、コントローラがオンライン除去されている間、冗長コントローラが静止状態になるか、あるいはリセット状態に保持されることによって生じる。

【0063】

冗長コントローラシステムからコントローラをオンライン除去するための知られているプロセスには早期警告スイッチすなわち短コネクタピンが含まれており、これはコントローラが除去されていることを冗長コントローラシステムに警告するものである。この警告により、コントローラはその時点のコントローラのメモリアクセスを終了させ、コントローラの不揮発メモリを自己リフレッシュモードにすることができる。

10

【0064】

コントローラが完全に分離（例えば、長コントローラ検出ピンが開放）されると、「対になっている」コントローラのもう一方のコントローラが、冗長コントローラシステムを制御するための動作を再開することができる。検出ピンコネクタを利用しているシステムの場合、部分除去されるコントローラは検出ピンが完全に接点をシャットダウンするまでの間、システムを使用不可の状態に保持することができる。このセットアップでは、例えばコントローラが冗長コントローラシステムから「部分的にのみ除去」される等の誤った手順により、オンライン除去故障時間が延長し、ホストオペレーティングシステムのタイムアウト期間を超過するという場合がある。

20

【0065】

図13は、本発明による冗長コントローラシステムからコントローラをオンライン除去する方法の一例示的实施形態を示したものである。上記方法は450で示されている。オンライン除去方法450により、冗長コントローラシステムの故障時間を最短にする一方で、コントローラを冗長コントローラシステムから安全にオンライン除去することができる。

【0066】

452で、冗長コントローラシステム110からコントローラが除去されていることが検出される。ここでも、本明細書において既に説明した図1から図12を参照する。一態様では、コネクタ上の早期警告スイッチ又は短コネクタピンが、コントローラが除去されていることをシステム110に警告するが、検出は冗長コントローラシステム110からコントローラが完全に開放される前になされることが好ましい。上記警告は、ホットプラグ警告142またはホットプラグ警告182を介して受信される。本明細書で説明する一例示的实施形態では、第1のコントローラ112が冗長コントローラシステム110から除去される。

30

【0067】

454で、冗長コントローラシステム110から第1のコントローラ112が除去されていることが検出されると、第1のコントローラ112および第2のコントローラ114に対するシャットダウンシーケンスが実行される。一態様では、456で、各コントローラに対するシャットダウンシーケンスでは、コントローラのプロセッサに割込みがかけられ、プロセッサによる稼働中の処理タスクの終了を可能にしている。458で、各コントローラに対するシャットダウンシーケンスではさらに、メモリ120およびメモリ160に対する未処理のメモリアクセスを完了し、内部バッファをフラッシングする。シャットダウンシーケンスの一部として、メモリコントローラ122によって状態語がメモリ120に書き込まれ、また、メモリコントローラ162によって状態語がメモリ160に書き込まれる。

40

【0068】

好ましい一実施形態では、第1のメモリ120および第2のメモリ160は、自己リフレッシュモードを有しており、さらに好ましくは、バッテリーバックアップを備えている。第1のコントローラ112および第2のコントローラ114に対するシャットダウンシーケ

50

ンスが完了すると、第1のメモリ120および第2のメモリ160が、それぞれ対応する第1のメモリコントローラ122および第2のメモリコントローラ162によって自己リフレッシュモードに置かれる。462で、除去を検出するコントローラはオフラインに留まり、除去の終了を待つ。そのメモリは、自己リフレッシュモードに留まっている。464で、除去を検出しないコントローラによって、オンラインになるプロセスが直ちに開始される。

【0069】

本明細書で説明する例示的实施形態では、メモリ120に対する自己リフレッシュプロセスが完了すると、除去を検出する第1のコントローラ112がオフラインに留まり、除去の終了を待ち、およびメモリ120が自己リフレッシュモードに留まる（一態様では、メモリは、バッテリーでバックアップされたDRAMである）。除去を検出しない第2のコントローラ114は、直ちにオンラインになるプロセスを開始し、冗長コントローラシステムの故障時間を最短にする。メモリ160（例えばバッテリーバックアップDRAM）は、自己リフレッシュモードを解除され、メモリコントローラが使用する、予め書き込まれている状態語を保持している。

【0070】

以上、好ましい実施形態について説明するために、本明細書において特定の実施形態を示し、かつ、記述したが、本発明の範囲を逸脱することなく、広範囲の様々な代替実施態様および/または等価実施態様を、上に示し、かつ、説明した特定の実施形態と置き換えることができることは、当分野の技術者には理解されよう。本発明を、極めて広範囲の様々な実施形態の中で実施することができることは、化学、機械、電気機械、電気、およびコンピュータ技術の分野の技術者には、容易に理解されよう。本出願は、本明細書において説明した好ましい実施形態のあらゆる応用、すなわち変形形態を包含することを意図している。したがって本発明は、特許請求の範囲の各クレームおよびその相当によってのみ制限されることを、明確に意図している。

【0071】

（1）第1のコントローラ（32、112）および第2のコントローラ（34、114）を有する冗長コントローラシステム（30、110）からコントローラをオンライン除去する方法であって、

早期警告検出を介して、前記冗長コントローラシステムからの前記第1のコントローラの部分除去を検出するステップ（452）と、

前記第1のコントローラおよび前記第2のコントローラに対して、未処理のメモリアクセスを完了させるステップ（458）を含むシャットダウンシーケンスを実行するステップ（454）と、

第1のメモリ（50、120）を有するように前記第1のコントローラを定義し、かつ、前記第1のメモリを自己リフレッシュモードにするステップ（460）と、前記冗長コントローラシステムからの前記第1のコントローラの除去を終了させるステップを含む方法。

【0072】

（2）第1のプロセッサ（54、124）を含むように前記第1のコントローラを定義するステップと、第2のプロセッサを含むように前記第2のコントローラを定義するステップとを含み、前記第1のコントローラおよび前記第2のコントローラに対してシャットダウンシーケンスを実行する前記ステップ（454）が、前記第1のプロセッサに割り込むステップおよび前記第2のプロセッサに割り込むステップ（456）をさらに含む、（1）に記載の方法。

【0073】

（3）前記第1のプロセッサに割り込む前記ステップが、前記第1のプロセッサによる未処理のプロセッサタスクの終了を可能にするステップを含み、かつ、前記第2のプロセッサに割り込む前記ステップが、前記第2のプロセッサによる未処理のプロセッサタスクの終了を可能にするステップ（456）を含む、（3）に記載の方法。

【 0 0 7 4 】

(4) 前記第 1 のメモリを前記自己リフレッシュモードにする前記ステップ (4 6 0) が、前記第 1 のメモリを、メモリの内容を維持するためのメモリコントローラからの外部リフレッシュサイクルを必要としないモードにするステップを含む、(1) に記載の方法。

【 0 0 7 5 】

(5) 第 2 のメモリを備えるように前記第 2 のコントローラを定義し、かつ、前記第 2 のメモリを自己リフレッシュモードにするステップ (4 6 0) を含む、(1) に記載の方法。

【 0 0 7 6 】

(6) 前記第 2 のメモリを自己リフレッシュプロセスに置く前記ステップ (4 6 0) が、前記第 2 のメモリを、メモリの内容を維持するためのメモリコントローラからの外部リフレッシュサイクルを必要としないモードにするステップを含む、(5) に記載の方法。

【 0 0 7 7 】

(7) 前記第 2 のメモリの自己リフレッシュモードを終了させた後で、前記第 2 のコントローラに、オンラインになるプロセスを開始させるステップ (4 6 4) を含む、(5) に記載の方法。

【 0 0 7 8 】

(8) 前記オンラインになるプロセスが、前記自己リフレッシュモードを終了させた後で、かつ、前記冗長コントローラシステムからの前記第 1 のコントローラの除去を終了させる前に、直ちに開始されるように構成される (4 6 4)、(7) に記載の方法。

【 0 0 7 9 】

(9) 第 1 のプロセッサ (5 4、1 2 4) を含むように前記第 1 のコントローラを定義し、かつ、第 2 のプロセッサ (6 0、1 6 4) を含むように前記第 2 のコントローラを定義するステップと、

前記第 1 のコントローラおよび前記第 2 のコントローラに対して、前記第 1 のプロセッサおよび前記第 2 のプロセッサに割り込み、前記第 1 のコントローラおよび前記第 2 のコントローラに対する未処理のメモリアクセスを完了させるステップ (4 5 6) を含むシャットダウンシーケンスを実行するステップ (4 5 4) と、

前記第 1 のメモリの自己リフレッシュモードを解除した後、前記冗長コントローラシステムからの前記第 1 のコントローラの除去を終了させるステップとを含む、(1) に記載の方法。

【 0 0 8 0 】

(1 0) (1) から (9) に記載の方法を実行するように構成された冗長コントローラシステム (3 0、1 1 0)。

【図面の簡単な説明】

【図 1】冗長コントローラをホット挿入するために構成された、本発明による冗長コントローラデータ格納システムの一例示的实施形態を示す図である。

【図 2】本発明による冗長コントローラデータ格納システムにコントローラをホット挿入するための方法の一例示的实施形態を示す図である。

【図 3】本発明による冗長コントローラシステムにコントローラをホット挿入する方法の他の例示的实施形態を示す図である。

【図 4】本発明による冗長コントローラをホット挿入するために構成された冗長コントローラデータ格納システムの他の例示的实施形態を示す構成図である。

【図 5】本発明による冗長コントローラデータ格納システムに使用されるタスクプロセッサの一例示的实施形態を示す図である。

【図 6】本発明による冗長コントローラシステムのタスクプロセッサによって利用されるデータ構造の一例示的实施形態を示す図である。

【図 7】本発明による冗長コントローラデータ格納システムに使用される、メモリブロック中に構成されたメモリイメージを有するコントローラ共用メモリの一例示的实施形態を示す図である。

10

20

30

40

50

【図 8】本発明による冗長コントローラシステムにコントローラをホット挿入する方法の一例示的实施形態を示す図である。

【図 9】本発明による冗長コントローラシステムにコントローラをホット挿入する方法の一例示的实施形態をさらに示す図である。

【図 10】本発明による冗長コントローラシステムにコントローラをホット挿入する方法の一例示的实施形態をさらに示す図である。

【図 11】本発明による冗長コントローラシステムにおけるコントローラリセット処理方法の一例示的实施形態を示す図である。

【図 12】本発明による冗長コントローラシステムにおけるコントローラリセット処理方法の一例示的实施形態をさらに示す図である。

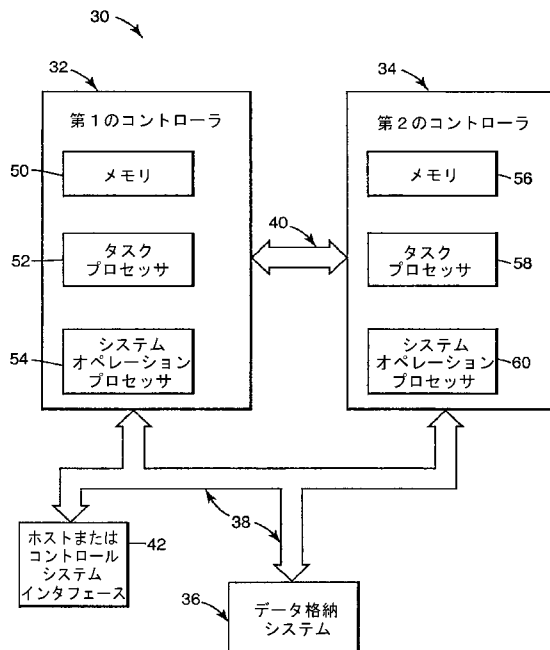
【図 13】本発明による冗長コントローラシステムからコントローラを除去する方法の一例示的实施形態を示す図である。

【符号の説明】

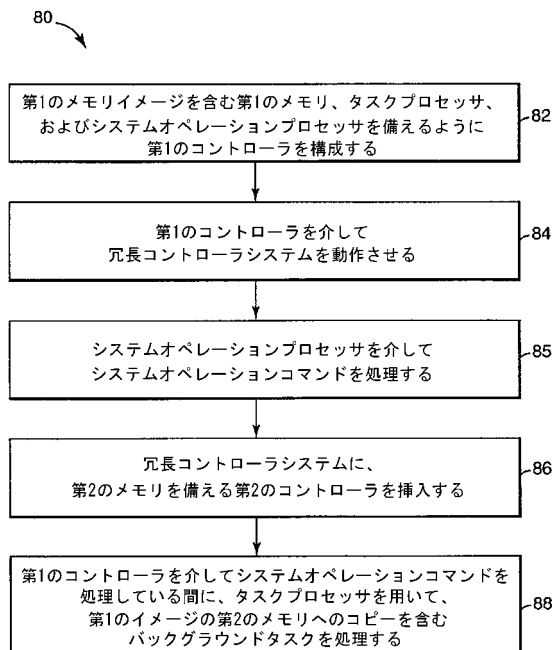
32、112	第1のコントローラ
30、110	冗長コントローラシステム
34、114	第2のコントローラ
50、120	第1のメモリ
54、124	第1のプロセッサ
60、164	第2のプロセッサ

10

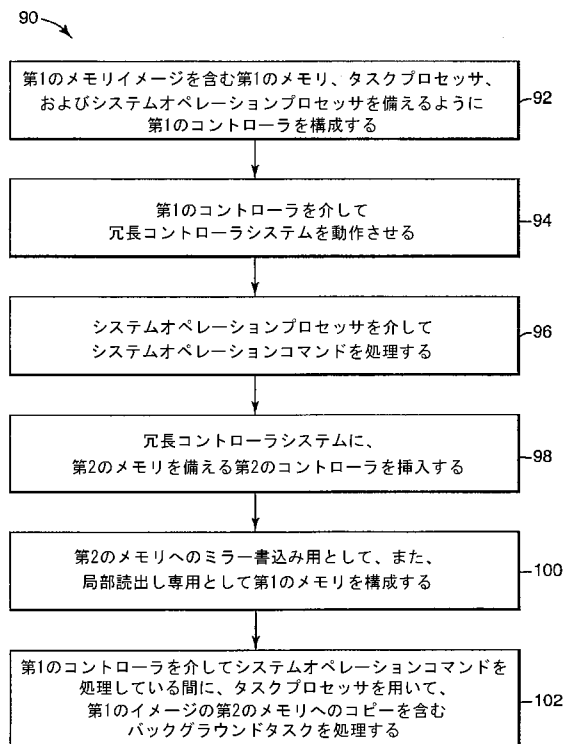
【図 1】



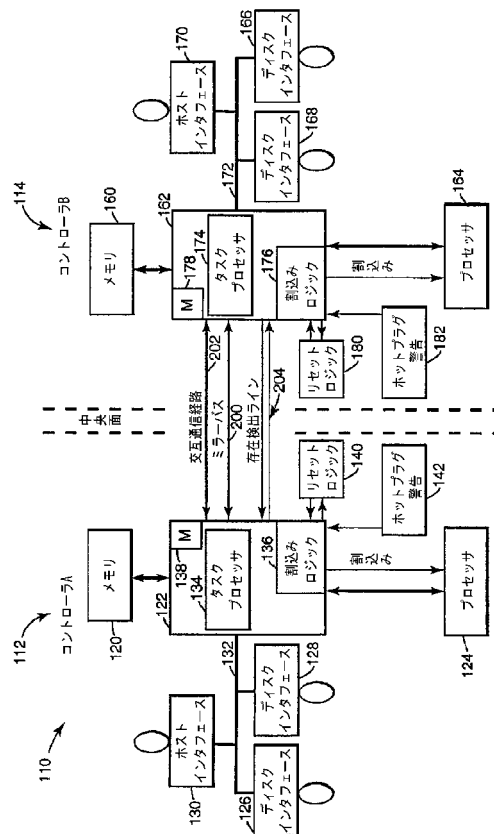
【図 2】



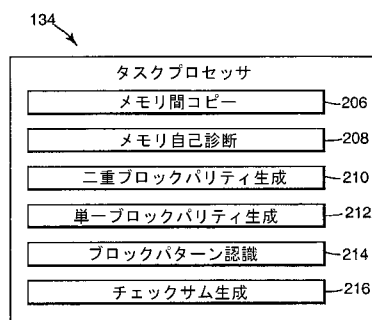
【図 3】



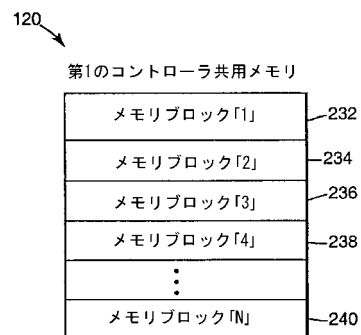
【図 4】



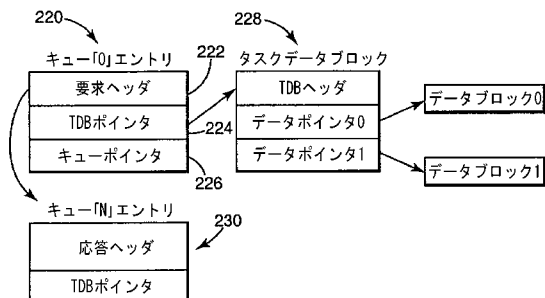
【図 5】



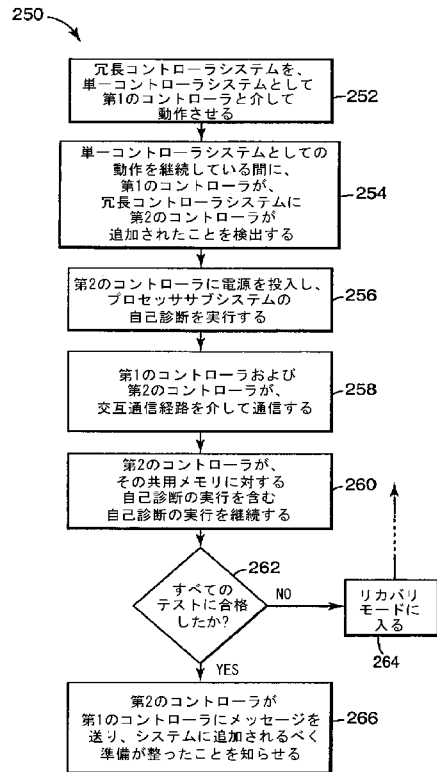
【図 7】



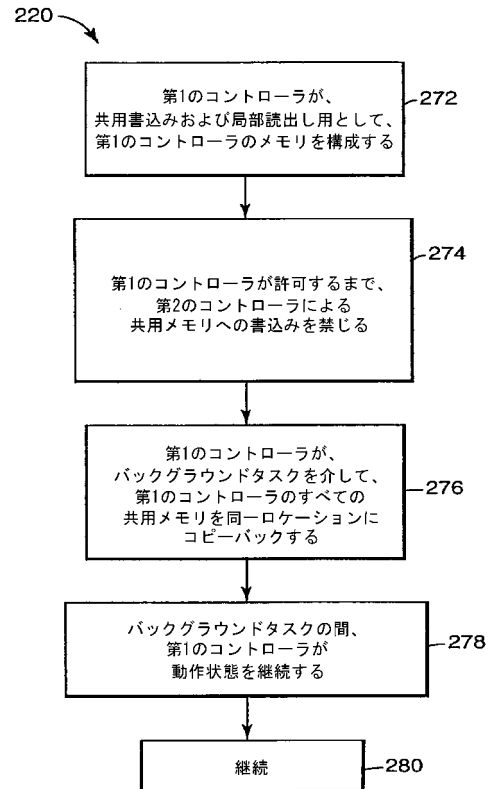
【図 6】



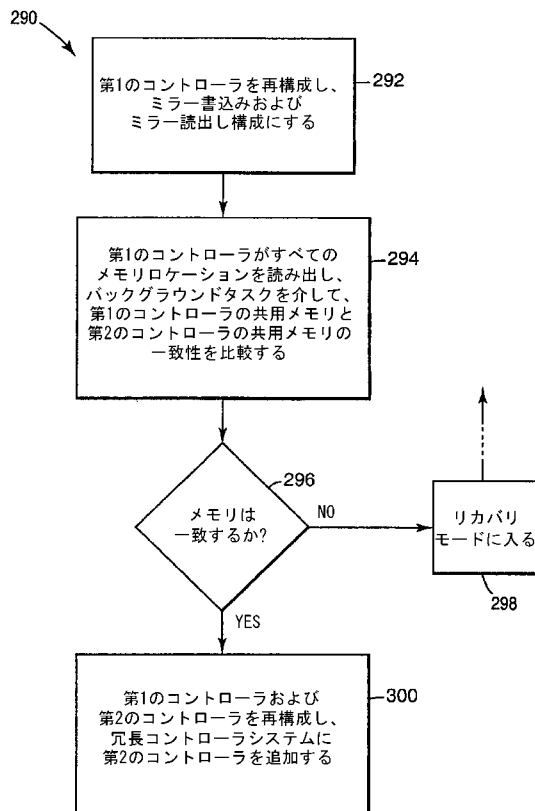
【図 8】



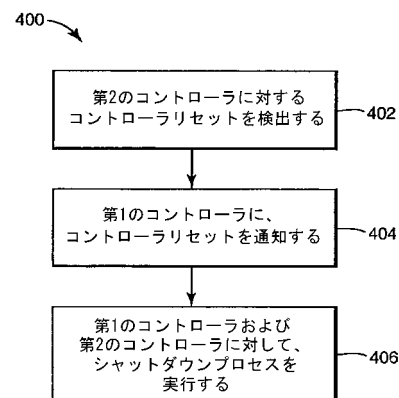
【図 9】



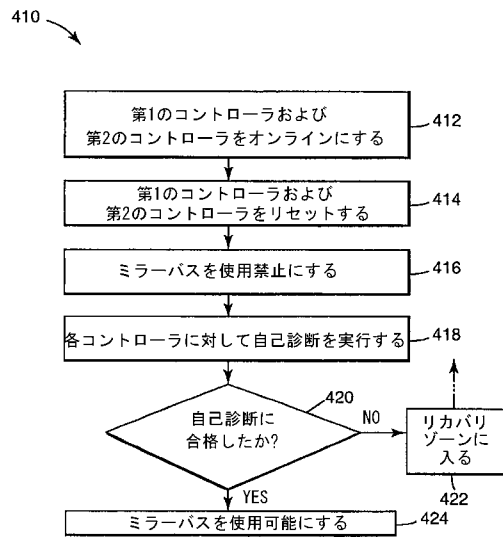
【図 10】



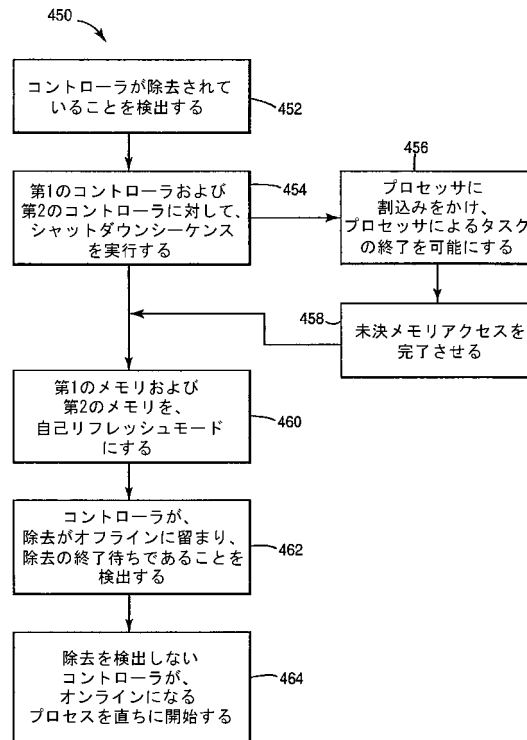
【図 11】



【図 12】



【図 13】



フロントページの続き

- (72)発明者 バリー・ジェイ・オールドフィールド
アメリカ合衆国 8 3 7 1 3 アイダホ州ボイジー、ウェスト・ダニエル・コート 1 1 3 0 2
- (72)発明者 クリストファー・ダブリュ・ヨハンソン
アメリカ合衆国 8 3 6 2 9 アイダホ州ホースシュー・ベンド、リバー・ブラフ・レーン 3 2

審査官 高 橋 正 徳

- (56)参考文献 特開平 0 8 - 2 7 2 7 5 3 (J P , A)
特開平 0 8 - 1 9 0 4 9 4 (J P , A)
特許第 2 6 0 8 9 0 4 (J P , B 2)
国際公開第 0 0 / 0 6 7 1 2 6 (W O , A 1)

- (58)調査した分野(Int.Cl. , D B 名)
G06F11/16-11/20