

(19) 日本国特許庁(JP)

(12) 公 開 特 許 公 報(A)

(11) 特許出願公開番号  
特開2004-21556  
(P2004-21556A)

(43) 公開日 平成16年1月22日 (2004.1.22)

(51) Int.Cl. <sup>7</sup>	F I	テーマコード (参考)
GO6F 11/16	GO6F 11/16 31OC	5B034
GO6F 3/06	GO6F 3/06 3O4N	5B042
GO6F 11/20	GO6F 11/20 31OE	5B065
GO6F 11/30	GO6F 11/30 31OH	

審査請求 未請求 請求項の数 21 O L (全 18 頁)

(21) 出願番号	特願2002-174944 (P2002-174944)	(71) 出願人	000005108
(22) 出願日	平成14年6月14日 (2002.6.14)		株式会社日立製作所
			東京都千代田区神田駿河台四丁目6番地
		(74) 代理人	100071283
			弁理士 一色 健輔
		(74) 代理人	100084906
			弁理士 原島 典孝
		(74) 代理人	100098523
			弁理士 黒川 恵
		(74) 代理人	100112748
			弁理士 吉田 浩二
		(74) 代理人	100110009
			弁理士 青木 康

最終頁に続く

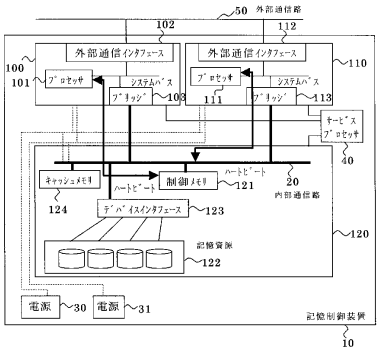
(54) 【発明の名称】 記憶制御装置およびその制御方法

(57) 【要約】 (修正有)

【課題】 内部通信路を用いて稼働情報を効率よく管理し、性能や信頼性が高く構築が容易で安価に移動情報の監視機能を実現する記憶制御装置を提供する。

【解決手段】 本発明の記憶制御装置は、外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段に対するデータ入出力指示を送信する複数の制御部と、前記制御部間のデータ入出力指示及びデータ入出力を行うための内部通信路と、前記制御部が前記内部通信路を通じて他の前記制御部に自身の稼働情報を送信することで前記制御部が互いに他の前記制御部の稼働状態を監視する手段と、を備える。

【選択図】 図1



## 【特許請求の範囲】

## 【請求項 1】

外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段に対するデータ入出力指示を送信する複数の制御部と、  
前記制御部間のデータ入出力指示及びデータ入出力を行うための内部通信路と、  
前記制御部が前記内部通信路を通じて他の前記制御部に自身の稼働情報を送信することで前記制御部が互いに他の前記制御部の稼働状態を監視する手段と、  
を備えることを特徴とする記憶制御装置。

## 【請求項 2】

外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段にデータ入出力指示を送信する複数の制御部と、  
前記制御部間におけるデータ入出力指示及びデータ入出力を行い前記記憶手段が接続する内部通信路と、  
前記制御部の稼働情報を前記制御部から前記内部通信路を通じて当該内部通信路に接続するメモリに送信しこれを当該メモリに記憶する手段と、  
前記制御部が前記内部通信路を通じて前記メモリにアクセスし前記稼働情報に基づいて他の前記制御部の稼働状況を監視する手段と、  
を備えることを特徴とする記憶制御装置。

## 【請求項 3】

前記制御部が、  
中央処理装置と、前記外部通信路に接続するための外部通信インタフェースと、  
前記内部通信路と接続するための内部通信インタフェースと、  
を備えることを特徴とする請求項 2 に記載の記憶制御装置。

## 【請求項 4】

前記稼働情報には前記制御部の障害有無を示す情報が含まれており、  
前記制御部が前記監視手段により他の前記制御部に障害が生じていることを認知した場合、障害が生じている前記制御部が行っていた処理を他の前記制御部に引き継ぐ手段を備えることを特徴とする請求項 1 または 2 のいずれかに記載の記憶制御装置。

## 【請求項 5】

前記メモリは前記制御部が担当する処理に関するリソース情報を記憶する手段を備え、  
前記処理の引き継ぎ先となる前記制御部が、前記リソース情報にアクセスし当該制御部自身が引き継ぐべき前記処理を認知する手段を備えることを特徴とする請求項 4 に記載の記憶制御装置。

## 【請求項 6】

前記リソース情報には、前記制御部が前記外部通信路もしくは前記内部通信路による通信に必要な情報が含まれることを特徴とする請求項 5 に記載の記憶制御装置。

## 【請求項 7】

前記リソース情報には、前記データ入出力指示において指定する記憶領域指定情報が含まれることを特徴とする請求項 5 に記載の記憶制御装置。

## 【請求項 8】

前記リソース情報には、外部通信路上における前記制御部のネットワークアドレスが含まれることを特徴とする請求項 5 に記載の記憶制御装置。

## 【請求項 9】

前記稼働情報には、前記中央処理装置、前記外部通信インタフェース、前記内部通信インタフェースのいずれかについての障害の有無を示す情報が含まれることを特徴とする請求項 3 に記載の記憶制御装置。

## 【請求項 10】

前記障害の有無を示す情報が、前記制御部が前記内部通信路を介して前記メモリに一定間隔で送信され、前記制御部に対応づけて記憶されるタイムスタンプであり、  
前記監視手段は、ある前記制御部に対応する前記タイムスタンプが一定時間以上更新され

ていない場合に、そのタイムスタンプに対応づけられている前記制御部に障害が生じていると認知することを特徴とする請求項 9 に記載の記憶制御装置。

【請求項 11】

前記監視手段が、前記稼働情報にアクセスできない場合に、前記内部通信路に障害が生じていると判断する手段を備えることを特徴とする請求項 4 に記載の記憶制御装置。

【請求項 12】

前記監視手段が、前記稼働情報にアクセスできかつ前記のある制御部の前記タイムスタンプが一定時間以上更新されていない場合に、その制御部自体に障害が生じていると判断する手段を備えることを特徴とする請求項 4 に記載の記憶制御装置。

【請求項 13】

前記記憶手段が、中央処理装置と、前記メモリと、ディスクドライブなどの記憶資源に対するデータ入出力を行うデバイスインタフェースとを備えることを特徴とする請求項 2 に記載の記憶制御装置。

【請求項 14】

前記稼働情報、もしくは、前記リソース情報は、前記メモリに記憶されることを特徴とする請求項 13 に記載の記憶制御装置。

【請求項 15】

前記稼働情報、もしくは、前記リソース情報は、前記記憶資源に記憶されることを特徴とする請求項 13 に記載の記憶制御装置。

【請求項 16】

前記制御部が、データ入出力要求をファイル名単位で行うファイルシステムを備え、ファイル名により指定されるデータを単位として前記記憶手段に対する前記データ入出力指示を送信する手段を備えることを特徴とする請求項 2 に記載の記憶制御装置。

【請求項 17】

前記記憶手段が、ハードディスクなどの記憶資源に対するデータ入出力を制御するデバイスインタフェースと、  
キャッシュメモリと、  
前記キャッシュメモリを介して前記記憶資源に対するデータの読み書きを行う手段と、  
を備えることを特徴とする請求項 2 に記載の記憶制御装置。

【請求項 18】

前記内部通信路が、複数の通信経路により冗長構成されることを特徴とする請求項 1 に記載の記憶制御装置。

【請求項 19】

前記制御部が、前記内部通信路とは別に設けられ前記制御部間を接続する他の通信路を通じて他の前記制御部に自身の稼働情報を送信する手段を備え、  
前記制御部が前記他の通信路を通じて送信しようとする前記稼働情報を前記内部通信路に対して送信する変換手段を備えることを特徴とする請求項 2 に記載の記憶制御装置。

【請求項 20】

外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段に対するデータ入出力指示を送信する複数の制御部と、前記制御部間のデータ入出力指示及びデータ入出力を行うための内部通信路と、を備える記憶制御装置の制御方法であって、  
前記制御部が前記内部通信路を通じて他の前記制御部に自身の稼働情報を送信することで前記制御部同士が他の前記制御部の稼働状態を監視することを特徴とする記憶制御装置の制御方法。

【請求項 21】

外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段にデータ入出力指示を送信する複数の制御部と、前記制御部間におけるデータ入出力指示及びデータ入出力を行い前記記憶手段が接続する内部通信路と、を備える記憶制御装置の制御方法であって、前記制御部の稼働情報を前記制御部から前記内部通信路を通じて当該内部通信路に接続するメモリに送信しこれを当該メモリに記憶し、

10

20

30

40

50

前記制御部が前記内部通信路を通じて前記メモリにアクセスし前記稼働情報に基づいて他の前記制御部の稼働状況を監視することを特徴とする記憶制御装置の制御方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

この発明は、外部通信路から入力されるデータ入出力指示を受信して記憶デバイスに対するデータ入出力制御を行い内部通信路で接続された複数の制御部を備える記憶制御装置に関し、とくに制御部間で交換する稼働情報の通信路として内部通信路を用い、稼働情報を効率よく管理し、性能や信頼性が高く構築が容易で安価に制御部間での稼働情報の監視機能を提供する技術に関する。

10

【0002】

【従来の技術】

近年、ストレージ製品のの一つとして、内部にファイルシステムが実装され、ファイル指定によるデータ入出力要求を取り扱うことができるようにした記憶制御装置である、いわゆる、NAS (Network Attached Storage) サーバに対するニーズが高まっている。また、このようなNASサーバとしては、処理能力の向上等を目的として、それぞれがLAN等の外部通信路に接続し、それぞれが個別に外部通信路を通じて送られてくる処理要求に応答できるようにした複数の制御部を、同一の筐体内に収容する構成の製品が存在する。

【0003】

20

【発明が解決しようとする課題】

ところで、記憶制御装置はミッションクリティカルな状況で使用される場合が多く、一般に高い可用性が要求される。

ここで、LAN上のコンピュータ同士の障害検知に際しては、従来から、LAN上のコンピュータ間でハートビートメッセージを交換して互いの稼働状態を監視し合うことで、可用性を向上させる仕組みが知られている。

【0004】

例えば、特開2000-222373号公報には、クラスタ化コンピュータシステムにおいて、クラスタを構成するコンピュータ同士でLANを通じて定期的にハートビートメッセージを交換し、お互いに稼働状態を監視し合うようにし、ハートビートメッセージの交換が正常に行われていない場合には、適宜制御部のデータサービスは他の正常な制御部へ引き継ぐようにした仕組みが開示されている。また、特開2001-100943号公報には、2台のPCサーバが、ハートビートを行うための通信を、ディスク装置が接続してあるSCSIバスを用いて行うようにしたクラスタシステムが開示されている。

30

【0005】

一方、前記図1の構成の記憶制御装置のように、複数の制御部を備える構成の記憶制御装置においては、各制御部は同一筐体内に実装されてはいるものの、外部通信路からの要求に対してそれぞれが個別にサービスを提供している。

【0006】

そこで、このような構成の記憶制御装置においても、前記公報のように、制御部間で稼働状態を監視し、また、ある制御部における障害を検知した場合には、その業務を他の正常な制御部に引き継ぐ仕組みを設けることは、可用性の向上に有効であると考えられる。

40

【0007】

ここでこのような仕組みを設ける場合、稼働情報を伝達する通信路に何を用いるかが問題となる。例えば、前記公報のクラスタ化コンピュータシステムの場合には、このような稼働情報(この場合はハートビートメッセージ)の通信路として、コンピュータ間を結ぶLAN等の外部通信路や、専用の通信路を用いている。

【0008】

しかしながら、外部通信路は、経路途中の信頼性が充分で無く、また、通信速度も充分でない、もしくは、一定しないといった問題がある。また、専用の通信路を設ける場合は、

50

専用の設備が必要となり、余分なコストが発生することになる。

【 0 0 0 9 】

一方、前述した図 1 の構成の記憶制御装置においては、制御部間は、例えば、回路基板上に形成され CPU とメモリを結ぶ制御バスのようにデータ入出力指示及びデータ入出力を行うための内部通信路により互いに接続されている。そして、内部通信路は一般に LAN などの外部通信路よりも伝送能力が高く信頼性も高い。従って、内部通信路を稼働情報の通信路として用いることで、信頼性の高い制御部間で稼働状態を監視する仕組みを、容易かつ安価に実現することは可能であると考えられる。

【 0 0 1 0 】

本発明は、このような観点に基づいてなされたもので、内部通信路を用いて稼働情報を効率よく管理し、性能や信頼性が高く構築が容易で安価に稼働情報の監視機能を実現する記憶制御装置を提供することを目的とする。 10

【 0 0 1 1 】

【課題を解決するための手段】

この目的を達成する本発明のうち主たる発明は、  
外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段に対するデータ入出力指示を送信する複数の制御部と、  
前記制御部間のデータ入出力指示及びデータ入出力を行うための内部通信路と、  
前記制御部が前記内部通信路を通じて他の前記制御部に自身の稼働情報を送信することで前記制御部が互いに他の前記制御部の稼働状態を監視する手段と、 20  
を備えることを特徴とする。

【 0 0 1 2 】

すなわち、一般に外部通信路に比べて信頼性が高く、高速大容量通信が可能な、例えば、回路基板上に形成され CPU とメモリを結ぶ制御バスのようにデータ入出力指示及びデータ入出力を行うための内部通信路を通じて稼働情報を通知するようにしたことで、記憶制御装置の部品点数の削減が図られ、装置コストを抑えて信頼性の高いシステムを実現できる。

【 0 0 1 3 】

【発明の実施の形態】

= = = 開示の概要 = = =

本明細書および添付図面の記載により、少なくとも、以下の事項が明らかとなる。  
外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段に対するデータ入出力指示を送信する複数の制御部と、  
前記制御部間のデータ入出力指示及びデータ入出力を行うための内部通信路と、  
前記制御部が前記内部通信路を通じて他の前記制御部に自身の稼働情報を送信することで前記制御部が互いに他の前記制御部の稼働状態を監視する手段と、  
を備えることを特徴とする記憶制御装置。 30

この記憶制御装置によれば、例えば、回路基板上に形成され CPU とメモリを結ぶ制御バスのようにデータ入出力指示及びデータ入出力を行うための内部通信路を通じて稼働情報（例えば、ハートビートメッセージ）を通知することで、記憶制御装置の部品点数の削減が図られ、装置コストを抑えて信頼性の高いシステムを実現できる。 40

【 0 0 1 4 】

外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段にデータ入出力指示を送信する複数の制御部と、前記制御部間におけるデータ入出力指示及びデータ入出力を行い前記記憶手段が接続する内部通信路と、前記制御部の稼働情報を前記制御部から前記内部通信路を通じて当該内部通信路に接続するメモリに送信しこれを当該メモリに記憶する手段と、前記制御部が前記内部通信路を通じて前記メモリにアクセスし前記稼働情報に基づいて他の前記制御部の稼働状況を監視する手段と、を備えることを特徴とする記憶制御装置。

このような構成の記憶制御装置によれば、制御部とは別体の前記メモリ（例えば、後述す 50

る制御メモリ)に存在する稼働情報(例えば、後述する「稼働状態管理テーブル」)に障害の状態を管理することが可能となるため、内部通信路の稼働状態に関する情報などと組み合わせて、障害の原因や発生部分をより細かく特定することができる。

【0015】

また、かかる記憶制御装置においては、例えば、前記制御部は、中央処理装置と、前記外部通信路に接続するための外部通信インタフェースと、前記内部通信路と接続するための内部通信インタフェースと、を備えることとする。

【0016】

また、かかる記憶制御装置においては、前記稼働情報には前記制御部の障害有無を示す情報が含まれており、前記制御部が前記監視手段により他の前記制御部に障害が生じていることを認知した場合、障害が生じている前記制御部が行っていた処理を他の前記制御部に引き継ぐ手段を備えることとする。これにより、記憶制御装置の可用性が確保されることになる。

10

【0017】

また、かかる記憶制御装置において、前記メモリは前記制御部が担当する処理に関するリソース情報を記憶する手段を備え、前記処理の引き継ぎ先となる前記制御部が、前記リソース情報にアクセスし当該制御部自身が引き継ぐべき前記処理を認知する手段を備えることとする。

【0018】

また、かかる記憶制御装置において、前記リソース情報には、例えば、前記制御部が前記外部通信路もしくは前記内部通信路による通信に必要な情報、前記データ入出力指示において指定する記憶領域指定情報、外部通信路上における前記制御部のネットワークアドレスが含まれることとする。

20

【0019】

また、かかる記憶制御装置において、前記稼働情報には、前記中央処理装置、前記外部通信インタフェース、前記内部通信インタフェースのいずれかについての障害の有無を示す情報が含まれることとする。

【0020】

また、かかる記憶制御装置において、前記障害の有無を示す情報が、前記制御部が前記内部通信路を介して前記メモリに一定間隔で送信され、前記制御部に対応づけて記憶されるタイムスタンプであり、前記監視手段は、ある前記制御部に対応する前記タイムスタンプが一定時間以上更新されていない場合に、そのタイムスタンプに対応づけられている前記制御部に障害が生じていると認知することとする。

30

【0021】

また、かかる記憶制御装置において、前記監視手段が、前記稼働情報にアクセスできない場合に、前記内部通信路に障害が生じていると判断する手段を備えることとする。また、前記監視手段が、前記稼働情報にアクセスできかつ前記のある制御部の前記タイムスタンプが一定時間以上更新されていない場合に、その制御部自体に障害が生じていると判断する手段を備えることとする。これにより、障害を細かく特定することが可能となる。

【0022】

また、かかる記憶制御装置において、前記記憶手段が、中央処理装置と、メモリと、ディスクドライブなどの記憶資源に対するデータ入出力を行うデバイスインタフェースとを備えることとする。

40

【0023】

また、前記稼働情報、もしくは、前記リソース情報は、前記メモリに記憶されることとする。

【0024】

また、かかる記憶制御装置において、前記稼働情報、もしくは、前記リソース情報は、前記記憶資源に記憶されることとしてもよい。

【0025】

50

また、かかる記憶制御装置において、前記制御部が、データ入出力要求をファイル名単位で行うファイルシステムを備え、ファイル名により指定されるデータを単位として前記記憶手段に対する前記データ入出力指示を送信する手段を備えることとする。すなわち、記憶制御装置が、例えば、NASサーバとして用いられる場合である。

【0026】

また、かかる記憶制御装置は、前記記憶手段が、ハードディスクなどの記憶資源に対するデータ入出力を制御するデバイスインタフェースと、キャッシュメモリと、キャッシュメモリを介して前記記憶資源に対するデータの読み書きを行う手段と、を備えることを特徴とする。また、かかる記憶制御装置は、前記内部通信路が、複数の通信経路により冗長構成されることとする。

10

【0027】

また、かかる記憶制御装置は、前記制御部が、前記内部通信路とは別に設けられ前記制御部間を接続する他の通信路（例えば、後述する専用通信路）を通じて他の前記制御部に自身の稼働情報を送信する手段を備え、前記他の通信路に対して送信される稼働情報を前記内部通信路に対して送信する変換手段（例えば、後述するエミュレーションドライバ）を備えることを特徴とする。これにより汎用のクラスタソフトウェアなどを用いて、簡単かつ安価に稼働状態監視の仕組みを実現することができる。

【0028】

また、本発明の記憶制御装置の制御方法は、外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段に対するデータ入出力指示を送信する複数の制御部と、前記制御部間のデータ入出力指示及びデータ入出力を行うための内部通信路とを備える記憶制御装置の制御方法であって、前記制御部が前記内部通信路を通じて他の前記制御部に自身の稼働情報を送信することで前記制御部同士が他の前記制御部の稼働状態を監視することをする。

20

【0029】

また、本発明の他の記憶制御装置の制御方法は、外部通信路を通じて入力されるデータ入出力要求を受信して記憶手段にデータ入出力指示を送信する複数の制御部と、前記制御部間におけるデータ入出力指示及びデータ入出力を行い前記記憶手段が接続する内部通信路と、を備える記憶制御装置の制御方法であって、前記制御部の稼働情報を前記制御部から前記内部通信路を通じて当該内部通信路に接続するメモリに送信しこれを当該メモリに記憶し、前記制御部が前記内部通信路を通じて前記メモリにアクセスし前記稼働情報に基づいて他の前記制御部の稼働状況を監視することとする。

30

【0030】

== 第1実施例 ==

< 装置構成 >

まず、本発明を前述の図1に示す記憶制御装置10に適用した実施例について説明する。

【0031】

記憶制御装置10は、例えば、回路基板上に形成されCPUとメモリを結ぶ制御バス（システムバス）のようにデータ入出力指示及びデータ入出力を行うための内部通信路20と、これに接続する複数の制御部100、110、記憶装置120、冗長構成された電源装置30、31を備える。制御部100、110や記憶装置120に接続するサービスプロセッサ40は、制御部100、110や記憶装置120の動作制御や各種設定、稼働状態監視などを行う。

40

【0032】

制御部100、110は、CPUなどで構成される中央処理装置としてのプロセッサ101、111、LANなどの外部通信路50に接続するための外部通信インタフェース102、112、内部通信路20に接続するブリッジなどで構成された内部通信インタフェース103、113などを備える。

【0033】

制御部100、110では、ファイルシステム（不図示）が稼働し、制御部100、11

50

0 は外部通信路 5 0 からファイル名指定によるデータ入出力要求を取り扱う。つまり、制御部 1 0 0 , 1 1 0 は、それぞれ LAN 上のファイルサーバとして機能するコンピュータとしての機能を備え、記憶制御装置 1 0 は前述の NAS サーバとして機能している。

#### 【 0 0 3 4 】

一方、記憶装置 1 2 0 は、システム管理情報などが記憶される制御メモリ 1 2 1、ハードディスクなどの記憶資源 1 2 2、制御部 1 0 0 , 1 1 0 から送信されてくる命令などに応じて記憶資源 1 2 2 に対するデータの書き込み / 読み出しを実行するデバイスインタフェース 1 2 3、キャッシュメモリ 1 2 4などを備える。

#### 【 0 0 3 5 】

なお、記憶資源 1 2 2 は、図 1 に示すように記憶制御装置 1 0 に内蔵されている場合もあるし、また、記憶制御装置 2 0 の外部の別筐体内に存在し、デバイスインタフェース 1 2 3 と適宜なインタフェースで接続していることもある。

#### 【 0 0 3 6 】

##### < 基本動作 >

記憶制御装置 1 0 の基本的な動作について説明する。

記憶制御装置 1 0 がホストコンピュータなどの外部装置（不図示）から外部通信路 5 0 を通じて入力されるデータ入出力要求を受信すると、プロセッサ 1 0 1 , 1 1 1 は、この要求に対応する指示コマンドやデータなどからなるデータ入出力指示を、内部通信路 2 0 を介して制御メモリ 1 2 1 に送信する。制御メモリ 1 2 1 はこれを受信して記憶する。

#### 【 0 0 3 7 】

ここで、例えば、前記データ入出力指示に含まれる前記指示コマンドがライト (Write) コマンドであった場合、デバイスインタフェース 1 2 3 は、内部通信路 2 0 を介してプロセッサ 1 0 1 , 1 1 1 にデータ送信要求を送信する。この要求を受信したプロセッサ 1 0 1 , 1 1 1 は、キャッシュメモリ 1 2 4 にライトデータを格納し、また、デバイスインタフェース 1 2 3 に対して割り込み要求を送信する。この割り込み要求を受信したデバイスインタフェース 1 2 3 は、適宜な機会に、キャッシュメモリ 1 2 4 上の前記ライトデータを記憶資源 1 2 2 に書き込む。

#### 【 0 0 3 8 】

一方、前記指示コマンドがリード (Read) コマンドであった場合、デバイスインタフェース 1 2 3 は、この指示コマンドに付帯指定される記憶資源 1 2 2 上の記憶領域に格納されているデータを読み出し、これをキャッシュメモリ 1 2 4 に格納し、さらに、読み出したデータを、内部通信路 2 0 を介してプロセッサ 1 0 1 , 1 1 1 に送信する。

#### 【 0 0 3 9 】

割り込み要求を受信したデバイスインタフェース 1 2 3 は、制御メモリ 1 2 1 に格納されている前記指示コマンドを参照し、キャッシュメモリ 1 2 4 内のデータと冗長データとを記憶資源 1 2 2 に転送する。

#### 【 0 0 4 0 】

##### < 監視機能 >

各制御部 1 0 0 , 1 1 0 は、互いに他の制御部 1 0 0 , 1 1 0 に障害が発生したかどうかを監視している。各制御部 1 0 0 , 1 1 0 がどの他のどの制御部 1 0 0 , 1 1 0 の監視を担当するかは、例えば、サービスプロセッサ 4 0 を介してオペレータなどが設定し、設定された情報は、制御メモリ 1 2 1 上に存在する図 2 に示す稼働状態管理テーブルの「監視対象の制御部」の項目 2 5 1 , 2 5 4 に登録される。

#### 【 0 0 4 1 】

各制御部 1 0 0 , 1 1 0 は、自身の監視対象として割り当てられている他の制御部 1 0 0 , 1 1 0 の稼働状態を監視する。この監視により自身が担当する制御部 1 0 0 , 1 1 0 に何らかの障害が生じていることを検知した場合、制御部 1 0 0 , 1 1 0 は、検知した障害の内容に応じた処理を実行する。

#### 【 0 0 4 2 】

以下、制御部 1 0 0 , 1 1 0 による監視機能と、異常を検知した場合に実行される記憶制

10

20

30

40

50



御装置 1 0 の機能について、制御部 1 1 0 の障害を制御部 1 0 0 が検知する場合を例として図 3 のフローチャートとともに説明する。

【 0 0 4 3 】

制御部 1 1 0 は、内部通信路 2 0 を通じて定期的（タイミングは任意に変更できる）に制御メモリ 1 2 1 にアクセスする（もしくは、ハートビートメッセージを送信する）。一方、記憶装置 1 2 0 は、前記アクセスがあると、アクセスのあった時刻をタイムスタンプとして稼働状態管理テーブルの制御部 1 1 0 のタイムスタンプの項目 2 5 5 に書き込む（S 3 1 1）。すなわち、制御部 1 1 0 に異常が無ければタイムスタンプは定期的に更新されることになる。

【 0 0 4 4 】

一方、制御部 1 0 0 は、稼働状態管理テーブルを参照するため、内部通信路 2 0 を通じて定期的（定期的以外にも任意に設定してもよい）に制御メモリ 1 2 1 にアクセスする（S 3 1 2）。

【 0 0 4 5 】

ここで制御部 1 0 0 は、制御メモリ 1 2 1 にアクセスできなかった場合（S 3 1 3）、内部通信路 2 0 に何らかの障害が発生していると判断し、また、内部通信路 2 0 に障害が発生している場合は、記憶装置 1 2 0 へのデータ入出力が正常に行えない状態にある可能性があるので、制御部 1 0 0 は、例えば、サービスプロセッサ 4 0 に指示を出すなどして、制御部 1 1 0 の記憶装置 1 2 0 に対するデータ入出力処理を停止させるとともに、この処理に関して制御部 1 1 0 が取得中のリソースを開放させる（S 3 1 4）。

【 0 0 4 6 】

ここで、リソースとは、例えば、制御部 1 0 0、1 1 0 が外部通信路 5 0 もしくは内部通信路 2 0 による通信に際し必要となる、ネットワークアドレス（例えば、IP アドレス）などの情報、制御部 1 0 0、1 1 0 が記憶手段に対して送信するデータ入出力指示において指定する記憶領域指定情報（例えば、制御部 1 1 0 がマウントしていた記憶領域に関する情報）などである。

【 0 0 4 7 】

一方、制御部 1 0 0 は、制御メモリ 1 2 1 にアクセスできた場合、稼働状態管理テーブルの制御部 1 1 0 のタイムスタンプ 2 5 5 を参照し（S 3 1 5）、アクセスした時刻とタイムスタンプとの差が一定時間以上であるかどうかを調べる（S 3 1 6）。

【 0 0 4 8 】

ここで差が一定時間以上の場合には、制御部 1 0 0 は、制御部 1 1 0 に障害が発生していると判断し、制御部 1 1 0 が担当している処理やリソースを制御部 1 0 0 が引き継ぐ（S 3 1 7）。

【 0 0 4 9 】

一方、差が一定時間に満たない場合には、さらに、稼働状態管理テーブル中の監視対象の制御部 1 1 0 の「状態」項目の内容 2 5 6 を参照する（S 3 1 8）。ここでその内容が『正常』である場合には、制御部 1 0 0 は、制御部 1 1 0 は正常に動作していると認知する（S 3 1 9、S 3 2 0）。他方、「状態」項目の内容 2 5 6 が『異常』であった場合には、制御部 1 1 0 に障害が発生したと判断し、前記と同様の方法により制御部 1 1 0 が担当している処理を引き継ぐ（S 3 1 9、S 3 1 7）。なお、稼働状態管理テーブルにおいて、『正常』もしくは『異常』は、ビット表現等の適宜な形式で記述される。

【 0 0 5 0 】

ところで、記憶制御装置 1 0 において、稼働状態管理テーブルにおける「状態」項目の内容は、つぎのように管理されている。

例えば、制御部 1 0 0、1 1 0 は、図 4 に示すように外部通信インタフェース 1 0 2 に障害が発生している場合、そのことを外部通信インタフェース 1 0 2、1 1 2 からプロセッサ 1 0 1、1 1 1 への直接の障害報告、もしくは、外部通信インタフェース 5 0 に出した処理命令がタイムアウトする、といったことで認知する。

【 0 0 5 1 】

10

20

30

40

50

制御部 1 0 0 , 1 1 0 は、外部通信インタフェース 1 0 2 , 1 1 2 に障害が発生していることを認知すると、内部通信路 2 0 を通じて制御メモリ 1 2 1 上の稼働状態管理テーブルにアクセスし、障害が発生している制御部 1 0 0 , 1 1 0 に対応する「状態」項目 2 5 3 , 2 5 6 に『異常』を書き込む。

【 0 0 5 2 】

一方、制御部 1 0 0 , 1 1 0 は、図 5 に示すように内部通信インタフェース 1 0 3 , 1 1 3 に障害が発生している場合、そのことを内部通信インタフェース 1 0 3 , 1 1 3 からプロセッサ 1 0 1 , 1 1 0 への障害報告、もしくは、内部通信インタフェース 1 0 3 , 1 1 3 に出した処理命令のタイムアウトなどにより認知する。

【 0 0 5 3 】

制御部 1 0 0 , 1 1 0 は、内部通信インタフェース 1 0 3 , 1 1 3 に障害が発生していることを認知すると、内部通信路 2 0 を通じて制御メモリ 1 2 1 上の稼働状態管理テーブルにアクセスし、障害が発生している制御部 1 0 0 , 1 1 0 に対応する「状態」項目に『異常』を書き込む。

【 0 0 5 4 】

図 6 はプロセッサ 1 0 1 に障害が発生している場合である。この場合、制御部 1 0 0 は稼働状態管理テーブルに状態を書き込むことができない。しかしながら、制御部 1 1 0 のプロセッサ 1 1 1 が制御メモリ 1 2 1 を参照するため、プロセッサ 1 0 1 に障害している場合でも、その障害を検知することができる。

【 0 0 5 5 】

以上の実施例では、制御部 1 0 0 が制御部 1 1 0 の監視を行う場合について説明したが、当然のことながら制御部 1 1 0 が制御部 1 0 0 を監視する場合も同様の処理により行われる。

【 0 0 5 6 】

以上の実施例においては、制御メモリ 1 2 1 の稼働状態管理テーブルへの書き込みや参照などの制御部 1 0 0 , 1 1 0 における障害監視のための通信を、内部通信路 2 0 を介して行っている。このため、従来の外部通信路 5 0 や専用の通信路による方式に比べ、障害監視のための通信を高速に行える。また、一般に内部通信路 2 0 は、LAN などの外部通信路 5 0 に比べ伝送能力や信頼性に優れるため、障害監視のための通信を迅速かつ確実に行うことができる。

【 0 0 5 7 】

また、本来、制御部 1 0 0 , 1 1 0 間やこれらと記憶装置 1 2 0 間で行われるデータ入出力処理等のために設けられている内部通信路 2 0 を、障害監視のための通信に流用しているので、専用の通信路を設ける場合のように余分なハードウェアを増設する必要が無く、障害監視のための通信の仕組みを、容易かつ安価に構築できる。

【 0 0 5 8 】

内部通信路 2 0 を介して制御部 1 0 0 , 1 1 0 に接続する、制御部 1 0 0 , 1 1 0 とは別体の、記憶装置 1 2 0 の制御メモリに存在する稼働状態管理テーブルに障害の状態を管理する仕組みであるため、内部通信路 2 0 の稼働状態に関する情報などと組み合わせて、障害の原因や発生部分をより細かく特定することができる。

【 0 0 5 9 】

また、制御部 1 0 0 , 1 1 0 が障害の発生を認知した場合における、処理引き継ぎのための各制御部 1 0 0 , 1 1 0 が取得しているリソースに関する情報についても制御メモリ 1 2 1 に管理されるため、制御部 1 0 0 , 1 1 0 に障害が発生した場合でも、記憶装置 1 2 0 は引き継ぐリソースを確認することができる。また、リソースの一元管理により管理負荷の軽減等が図られる。

【 0 0 6 0 】

また、内部通信路が冗長構成されている場合には、障害監視のための通信の安全性や確実性がさらに担保される。

【 0 0 6 1 】

10

20

30

40

50

なお、以上の実施例は、制御部が２つの場合であったが、これに限定されるものではなく、制御部が３つ以上の場合にも容易に拡張することができる。

【００６２】

＝ ＝ 既存クラスタシステムの利用 ＝ ＝

前述の公報にも記載されているように、従来から、外部通信路や専用の通信路を利用してコンピュータ間で障害監視のための通信を行うソフトウェアが存在する。ここでは、このようなソフトウェアを用いて障害監視のための通信に、内部通信路２０を利用する本発明の仕組みを実現する場合について説明する。

【００６３】

図７は、本発明の制御部に対応する、コンピュータ７２０，７３０が、専用の通信路により接続され、この専用通信路７５０を介してハートビートメッセージを伝送することで、コンピュータ７２０，７３０がお互いに稼働状態を監視し合う、従来のシステム構成である。 10

【００６４】

コンピュータ７２０，７３０は、それぞれ、クライアントからのファイルサービスの要求をネットワークインタフェース７２１、７３１、ネットワークドライバ７２２，７３２を経由しファイルシステム７２３，７３３で受信する。

【００６５】

ファイルシステム７２３，７３３は、記憶装置７２６に対するデータ転送が必要な場合、ストレージドライバ７２４，７３４、ストレージインタフェース７２５，７３５を経由し 20  
記憶装置７２６とデータ転送を行う。クラスタソフトウェア７２７，７３７は、ハートビート用ネットワークドライバ７２８，７３８、ハートビートネットワークインタフェース７２９，７３９を経由し、コンピュータ７３０上のクラスタソフトウェア７２７，７３７に対しハートビートメッセージを送信する。

【００６６】

このクラスタソフト７２７，７３７が導入されている図７の構成に、障害監視に内部通信路を利用する本発明を適用した場合が図８である。この図において、記憶制御装置８１０は、クラスタソフトウェア８４７，８５７からのアクセスを受領し、このアクセスをストレージドライバ８４４，８５４へのアクセスに変換するエミュレーションドライバ８４８，８５８を備えている。この構成により、汎用のクラスタソフトウェア８４７，８５７が 30  
発行したハートビートメッセージを内部通信路８４９を通じて伝送される稼働状態情報に変換すること、および、その逆の変換をすることができる。具体的には、例えば、エミュレーションドライバ８４８，８５８は、ＲＳ－２３２ＣやＬＡＮなどの他の通信路の伝送手順に従った通信によりクラスタソフトウェア８４７，８５７からのアクセスを受領して、これを内部通信路８４９上での通信に変換したり、逆に内部通信路８４９上での通信を前記他の通信路の伝送手順に変換してクラスタソフトウェア８４７，８５７に伝えるといった役割を果たす。なお、ＲＳ－２３２ＣやＬＡＮなどの他の通信路は、記憶制御装置１０がハードウェア／ソフトウェアとして実際に備えていてもよいし、備えていなくてもよい。また、前記変換機能部分は、エミュレーションドライバ８４８，８５８の内部に設けるのでは無く、ストレージドライバ８４４，８５４や記憶装置８４６における内部通信イン 40  
タフェース（不図示）の機能を提供するファームウェア（不図示）などに設けるようにしてもよい。

【００６７】

以上に説明したように、既存のクラスタソフトウェアが導入されているシステムにおいては、エミュレーションドライバを導入するだけで本発明を実施することができる。

【００６８】

また、既存のクラスタソフトウェアが導入されていない場合には、既存のクラスタソフトウェアにエミュレーションドライバを組み合わせて導入することで、クラスタソフトウェアの機能部分に対する開発費が抑えられるため、低廉なコストで本発明を実施することができる。 50

## 【 0 0 6 9 】

＝ ＝ 内部通信路の他の構成 ＝ ＝

内部通信路が制御バスで無く、ファイバチャネル ( F i b r e C h a n n e l )、I n t e l l i g e n t I / O、R a p i d I / Oなどの他のプロトコルに準拠した通信路で構成される場合を説明する。

図 9 は内部通信路をファイバチャネル ( F i b r e C h a n n e l ) プロトコルによる通信路で構成した記憶制御装置 9 7 1 の構成を示す図である。制御部 9 5 0、9 6 0 は、内部通信インタフェースとしてファイバチャネルインタフェース 9 5 3、9 6 3 を備える。

## 【 0 0 7 0 】

一方、記憶装置 9 8 0 もファイバチャネルインタフェース 9 8 6、9 9 6 を備え、制御部 9 5 0、9 6 0 とはファイバチャネルスイッチ 9 5 6、9 6 6、ファイバチャネルインタフェース 9 8 6、9 9 2 を介して接続する。

10

## 【 0 0 7 1 】

制御部 9 5 0、9 6 0 は、ファイバチャネルスイッチ 9 5 6、9 6 6、ファイバチャネルインタフェース 9 5 3、9 6 3 を介して制御部 9 5 0、9 6 0 に接続され、制御メモリ 9 8 2、9 9 2、書き込みデータ及びディスクドライブからの読出しデータを一時バッファリングするキャッシュメモリ 9 8 3、9 9 3 を備える。

## 【 0 0 7 2 】

制御部 9 5 0、9 6 0 は、デバイスインタフェース 9 9 4 を介してディスクドライブ群 9 8 5 へ接続される。制御部 9 5 0、9 6 0、記憶装置 9 8 0 は、同一の筐体の実装される冗長構成の電源 9 7 2、9 7 3 より給電され、保守機構 9 7 4 により、各種動作設定や稼働管理が行われる。

20

## 【 0 0 7 3 】

制御部 9 5 0 がディスクドライブ群からなる記憶資源 9 8 5 とライト ( W r i t e ) データ転送を行う場合を例として、制御の流れ、データの流れを説明する。制御部 9 5 0 がサービスネットワーク 9 7 0 によりデータサービスの依頼を受けると、プロセッサ 9 5 2 は、サービス依頼を記憶装置 9 8 0 に対する I / O コマンドに変換し、ファイバチャネルインタフェース 9 5 3、ファイバチャネルスイッチ 9 5 6、ファイバチャネルインタフェース 9 8 6 を介して、記憶装置 9 8 0 へ I / O コマンドを送信する。

## 【 0 0 7 4 】

送信された I / O コマンドは、制御メモリ 9 8 2 に格納される。I / O コマンドがライトコマンドである場合、ファイバチャネルインタフェース 9 8 6 は、プロセッサ 9 5 2 にデータ転送を指示する。この指示によりプロセッサ 9 5 2 は、ファイバチャネルプロトコルに従い、記憶装置 9 8 0 に書き込みデータを送信する。記憶装置 9 8 0 は、送信されてきた書き込みデータを受信して、これを一旦キャッシュメモリ 9 8 3 に格納する。この書き込みデータは、デバイスインタフェース 9 8 4 により冗長データと共に記憶資源 9 8 5 に転送される。

30

## 【 0 0 7 5 】

つぎに、図 9 のシステムにおいて、制御部 9 5 0 が障害になった場合には制御部 9 6 0 で、制御部 9 6 0 が障害になった場合には制御部 9 5 0 で、それぞれ障害が発生している制御部のリソースを引継ぎ、自動的にデータサービスなどの記憶制御装置 1 0 の業務を続行できるようにした、記憶制御装置 9 7 1 の仕組みについて説明する。

40

## 【 0 0 7 6 】

記憶資源 9 8 5 には、図 1 0 に示す稼働状態管理テーブルが格納されている。稼働状態管理テーブルには、各制御部を特定する情報 1 0 0 1、1 0 0 2、監視すべき対象の制御部を特定する情報 1 0 0 3、1 0 0 4 についての項目がある。これらの項目は、例えば、フェイルオーバーポリシーに従って設定される。タイムスタンプ 1 0 0 5、1 0 0 6 は、制御部 9 5 0、9 6 0 が稼働情報 ( ハートビートメッセージ ) の I / O を発行する際の制御部上の時刻を示す。また、この稼働状態管理テーブルの「状態」1 0 0 7、1 0 0 8 項目には、制御部 9 5 0、9 6 0 の『正常』、『異常』を示す情報がセットされている。これら

50

の項目にセットされる情報には、必要に応じてクラスタシステムを構築する際のサービスネットワーク 1070 上の制御部に割り当てられた名称などの固有の識別子が用いられる。

制御部 950 上のプロセッサ 952 は、適宜なタイミング（例えば、定期的に）で記憶資源 985 の稼働情報管理テーブルの読み出し命令を送信し、稼働情報管理テーブルに示す全データもしくは一部のデータを取得する。稼働情報管理テーブルを読み出すと、プロセッサ 952 は、全データから制御部ネーム領域 1001, 1002 を参照し、制御部 950 が有する制御部ネーム「Server A」に対応づけられている監視対象制御部の制御部ネームを参照する。そして、この例では、「Server A」の監視担当制御部として「Server B」が割り当てられているので「Server B」のタイムスタンプ 1006 と状態情報 1008 とを参照し、フェイルオーバー処理を実行する必要性の有無を判断する。

10

#### 【0077】

プロセッサ 952 は、「Server B」のタイムスタンプ 1006 と現在時刻との差が一定時間以上であった場合、もしくは、「状態」項目 1008 に『異常』を示す情報がセットされていた場合には、制御部 960 に何らかの障害が発生していると認知する。また、プロセッサ 952 は、稼働状態管理テーブルにアクセスできることで、内部通信路に障害が発生していないと認知し、障害が制御部 960 におけるものであると判断して制御部 960 が行っていた処理やリソースの引き継ぎを開始する。

#### 【0078】

なお、プロセッサ 952 は、制御部 960 の障害を検出するかどうかに関わらず、「Server A」のタイムスタンプと状態情報である図 10 の 1005, 1007 の項目の内容を適宜なタイミング（例えば定期的に）更新する。

20

#### 【0079】

制御部 950, 960 が定期的に以上の処理を行うことにより汎用 I/O インタフェースを通信路として用いた障害監視等の仕組みが実現されることになる。

#### 【0080】

つぎに、ファイバチャネルスイッチ 956 に障害が生じた場合についての記憶制御装置の動作について説明する。

ファイバチャネルスイッチ 956 が障害となった場合、制御部 950 は I/O 入出力が不可となるが、通常のフェイルオーバー方式により制御部 960 の処理を継続させた場合には制御部 960 に処理が集中し、性能が劣化する可能性がある。

30

#### 【0081】

そこで、このような場合には、汎用 I/O 切り替えソフトウェア等により、制御部 950 の入出力経路を、ファイバチャネルインタフェース 953 ファイバチャネルスイッチ 956 ファイバチャネルインタフェース 986 という通常の経路から、ファイバチャネルインタフェース 954 ファイバチャネルスイッチ 966 ファイバチャネルインタフェース 996 という経路に変更して I/O を継続する。また、同時に、ディスクボリューム群 985 に制御部 950 から行われていた制御部 950 のアクセス状態情報を格納、確認するための I/O も経路を切り替える。

40

#### 【0082】

この例では、制御部 950 は、ファイバチャネルインタフェース 954 ファイバチャネルスイッチ 966 ファイバチャネルインタフェース 996 を順に経由して、図 10 におけるタイムスタンプ 1005、制御部 950 の「状態」項目の内容 1007 をディスクボリューム群 985 に格納する。

#### 【0083】

なお、以上はファイバチャネルスイッチ 956 が障害となった場合を例として説明したが、ファイバチャネルインタフェース 953、ファイバチャネルスイッチ 956、ファイバチャネルインタフェース 986 等に障害が発生した場合においても以上の処理を拡張できる。

50

## 【 0 0 8 4 】

以上のように、この実施例では、制御部 9 5 0 , 9 6 0 の状態情報を搭載した稼働情報の伝送を、外部通信路や専用の通信路では無く、記憶装置 9 8 0 への汎用 I / O インタフェース上で行っている。これによりハートビート用ネットワークインタフェースの導入が必要でなく、その分、手間やコストが削減され、また、信頼性の高い汎用 I / O インタフェースの利用により、信頼性の高い稼働状態監視機能を実現できる。

## 【 0 0 8 5 】

また、制御部の稼働状態を格納するディスクドライブ上に、クラスタ動作に必要な制御部ネットワークアドレス等のネットワーク属性、データが格納されているディスクドライブ等のリソース引継ぎ情報を格納することにより、各制御部より読み出し可能となるため、リソース引継ぎの設定管理を一元化し、管理コストの削減を図ることができる。 10

## 【 0 0 8 6 】

また、図 9 では 2 つのファイバチャネルスイッチ 9 5 6 , 9 6 6 を含む構成であったが、ファイバチャネルスイッチを全く含まない構成や、3 つ以上のファイバチャネルスイッチを含む構成である場合にも本発明を拡張して適用することができる。また、ファイバチャネル以外の汎用 I / O インタフェースを用いる場合にも適用できる。

## 【 0 0 8 7 】

＝ ＝ ＝ 稼働情報の格納方式 ＝ ＝ ＝

つぎに、汎用クラスタソフトを使用して本発明のシステムを構築する際に、汎用クラスタソフトが送信する稼働情報を、記憶装置のディスクドライブなどの記憶資源（例えば、ディスクドライブ）上に格納するようにした実施例について説明する。 20

## 【 0 0 8 8 】

図 1 1 はこのような構成の記憶制御装置 1 0 0 0 の一例である。制御部 1 1 1 0 , 1 1 5 0 は、クライアントからのファイルサービスの依頼をネットワークインタフェース 1 1 1 1 , 1 1 5 1、ネットワークドライバ 1 1 1 2 , 1 1 5 2 を経由してファイルシステム 1 1 1 3 , 1 1 5 3 で受信する。

## 【 0 0 8 9 】

ファイルシステム 1 1 1 3 , 1 1 5 3 は、ストレージドライバ 1 1 1 4 , 1 1 5 4、ファイバチャネルインタフェース 1 1 1 5 , 1 1 5 5、ファイバチャネルスイッチ 1 1 1 6 , 1 1 5 6 を経由して記憶装置 1 1 1 9 とデータ転送を行う。ストレージドライバ 1 1 1 4 , 1 1 5 4 内には、汎用 I / O 切り替えソフトウェアが組み込まれている。 30

## 【 0 0 9 0 】

データの転送経路に障害が発生した場合、ストレージドライバ 1 1 1 4 , 1 1 5 4 は、データ転送の経路をファイバチャネルインタフェース 1 1 1 7 , 1 1 5 7、ファイバチャネルスイッチ 1 1 1 6 , 1 1 5 6 を経由する経路に切り替える。

## 【 0 0 9 1 】

このシステムにおいて、汎用クラスタソフトウェアは、ネットワークに対し制御部の状態情報を含む稼働情報を送信する。ここで当該稼働情報の要求をネットワークドライバとして受信し、汎用 I / O インタフェース経由で記憶装置へデータ転送を行う要求に変換するエミュレーションドライバ 1 1 2 0 , 1 1 6 0 を、各制御部において稼働させることで、汎用クラスタソフトウェアに対し、透過的に汎用 I / O インタフェース上でのハートビートの交換を行えるようにしている。 40

## 【 0 0 9 2 】

また、ストレージドライバ 1 0 1 4 は、汎用 I / O インタフェース上のハートビートである I / O に対しても、データ転送の経路障害時に経路をファイバチャネルインタフェース 1 1 1 7 , 1 1 5 7、ファイバチャネルスイッチ 1 1 1 6 , 1 1 5 6 を経由する経路に切り替えてデータ転送を継続する。これにより信頼性の高い稼働状態監視のための仕組みが実現される。

## 【 0 0 9 3 】

## 【 発明の効果 】

本発明の記憶制御装置にあっては、制御部間で交換する稼働情報の通信路として内部通信路を用い、稼働情報を効率よく管理し、性能や信頼性が高く構築が容易で安価に稼働情報の監視機能を提供することができる。

【図面の簡単な説明】

【図１】本発明の一実施例による記憶制御装置の構成を示す図である。

【図２】本発明の一実施例による稼働情報管理テーブルを示す図である。

【図３】本発明の一実施例として説明する記憶制御装置における、制御部が他の制御部が検知する処理を説明するフローチャートである。

【図４】本発明の一実施例として説明する記憶制御装置における、外部通信インタフェースに障害が発生している場合の障害認知等の処理を説明する図である。

10

【図５】本発明の一実施例として説明する記憶制御装置における、内部通信インタフェースに障害が発生している場合の障害認知等の処理を説明する図である。

【図６】本発明の一実施例として説明する記憶制御装置における、プロセッサに障害が発生している場合の障害認知等の処理を説明する図である。

【図７】従来における、コンピュータがお互いに稼働状態を監視し合う仕組みを説明する図である。

【図８】本発明の一実施例による、障害監視に内部通信路を用いた記憶制御装置の構成を示す図である。

【図９】本発明の一実施例による、内部通信路をファイバチャネルプロトコルの通信路で構成した記憶制御装置の構成を示す図である。

20

【図１０】本発明の一実施例による稼働情報管理テーブルを示す図である。

【図１１】本発明の一実施例による、汎用クラスタソフトを使用し、また、稼働情報をディスクドライブなどの記憶資源に格納するようにした記憶制御装置の構成を示す図である。

【符号の説明】

２０ 内部通信路

５０ 外部通信路

１０１，１１０ 制御部

１２１ 制御メモリ

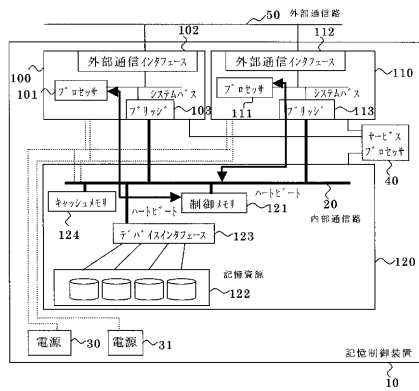
１２４ キャッシュメモリ

１２２ 記憶資源

１２３ デバイスインタフェース

30

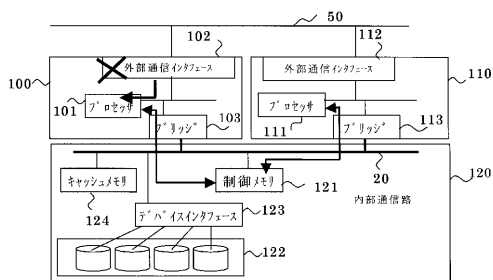
【図 1】



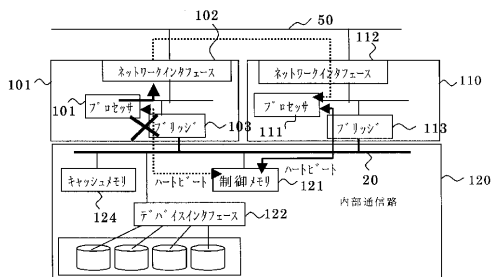
【図 2】

251			
制御部	監視対象の制御部	タイムスタンプ	状態
100	110	2001/11/13 10:05:02	正常
110	100	2001/11/13 10:02:15	正常

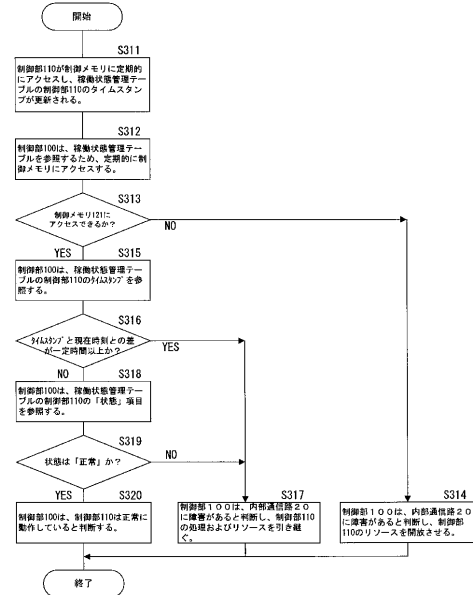
【図 4】



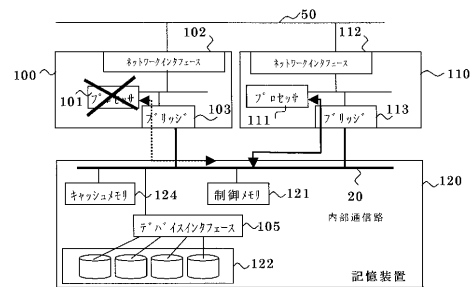
【図 5】



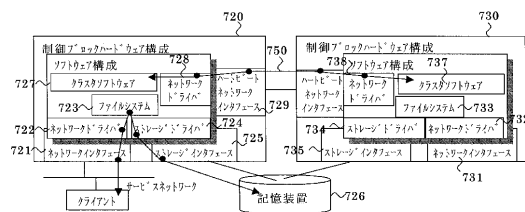
【図 3】



【図 6】

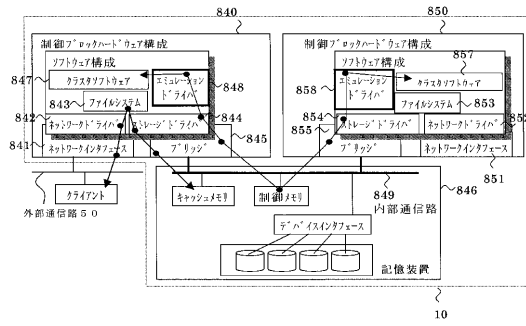


【図 7】



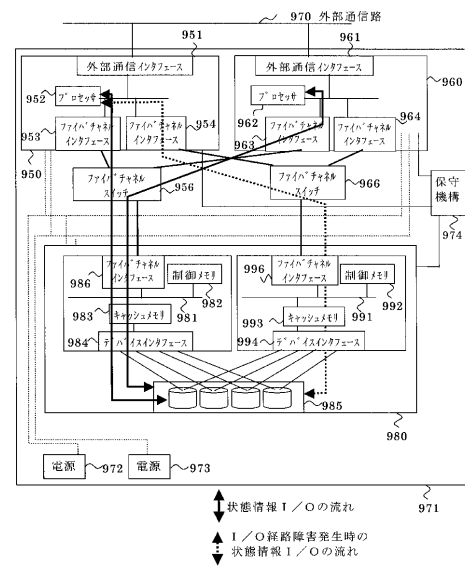


【図 8】



【図 9】

クラスシステムのハードウェア構成図



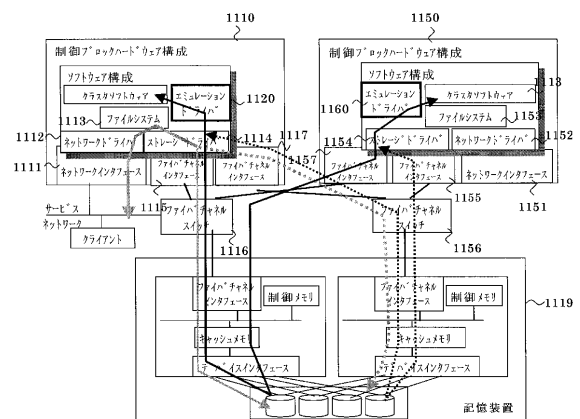
【図 10】

記憶装置内ハートビート用データ例

制御部ネーム	監視対象制御部	タイムスタンプ	状態
ServerA	ServerB	2001/11/13 10:05:02	正常
ServerB	ServerA	2001/11/13 10:02:15	正常

1001 1003 1005 1008  
1002 1004 1006 1007

【図 11】



---

フロントページの続き

(72)発明者 室谷 暁

神奈川県小田原市中里 3 2 2 番地 2 号 株式会社日立製作所 R A I D システム事業部内

(72)発明者 中野 俊夫

神奈川県小田原市中里 3 2 2 番地 2 号 株式会社日立製作所 R A I D システム事業部内

(72)発明者 横畑 静生

神奈川県小田原市中里 3 2 2 番地 2 号 株式会社日立製作所 R A I D システム事業部内

(72)発明者 高 本 賢一

神奈川県小田原市中里 3 2 2 番地 2 号 株式会社日立製作所 R A I D システム事業部内

F ターム(参考) 5B034 BB01 CC05 DD02

5B042 GA11 JJ04

5B065 BA01 CA11 EK02