

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5561622号
(P5561622)

(45) 発行日 平成26年7月30日(2014.7.30)

(24) 登録日 平成26年6月20日(2014.6.20)

(51) Int.Cl.		F I			
G06F	11/16	(2006.01)	G06F	11/16	310C
G06F	11/30	(2006.01)	G06F	11/30	305A
G06F	1/26	(2006.01)	G06F	1/00	330C
G06F	1/28	(2006.01)	G06F	1/00	333C

請求項の数 11 (全 22 頁)

(21) 出願番号	特願2011-210049 (P2011-210049)	(73) 特許権者	000004237
(22) 出願日	平成23年9月27日 (2011.9.27)		日本電気株式会社
(65) 公開番号	特開2013-73289 (P2013-73289A)		東京都港区芝五丁目7番1号
(43) 公開日	平成25年4月22日 (2013.4.22)	(74) 代理人	100102864
審査請求日	平成25年2月6日 (2013.2.6)		弁理士 工藤 実
		(72) 発明者	馬場 紀圭
			東京都港区芝五丁目7番1号 日本電気株式会社内
		審査官	▲高▼橋 正▲徳▼

最終頁に続く

(54) 【発明の名称】 多重化システム、データ通信カード、状態異常検出方法、及びプログラム

(57) 【特許請求の範囲】

【請求項1】

複数の物理マシンと、
前記複数の物理マシンの各々に搭載され、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行うためのデータ通信カードと
を含み、
前記データ通信カードは、
メイン電源の給電により駆動し、前記自系の物理マシン及び前記他系の物理マシンの内部状態を監視し、状態異常を検出する第1のLSI
を具備し、
前記第1のLSIは、
前記メイン電源がONの状態、監視及び検出対象の物理マシンのメモリに対してデータの読み出し要求を定期的に発行し、前記読み出し要求に対する応答を確認することで、前記監視及び検出対象の物理マシンの状態異常を検出する
多重化システム。

【請求項2】

請求項1に記載の多重化システムであって、
前記データ通信カードは、
搭載された物理マシンの基板上にチップとして搭載され、前記基板上にそれぞれ配置されているCPUとI/Oチップとの間に設けられており、

前記第1のLSIは、

前記メイン電源がONの状態、監視及び検出対象の物理マシンの前記I/Oチップに対して設定情報の読み出し要求を定期的に発行し、前記読み出し要求に対する応答を確認することで、前記監視及び検出対象の物理マシンの状態異常を検出する

多重化システム。

【請求項3】

請求項1または2に記載の多重化システムであって、

前記データ通信カードは、

スタンバイ電源の給電により駆動し、前記メイン電源がOFFの状態の時でも、前記自系の物理マシン及び前記他系の物理マシンの電源状態を監視し、電源の異常を検出する第2のLSI

を更に具備する

多重化システム。

【請求項4】

請求項3に記載の多重化システムであって、

前記第1のLSIは、

前記データ通信カード内部を制御するプロセッサと、

前記自系の物理マシン上で動作するソフトウェア(SW)の状態を取得するSW状態取得部と、

PCI Expressバスを介して、前記自系の物理マシン内部のハードウェア及びI/Oチップと接続し、前記ハードウェア及び前記I/Oチップの状態を監視するPCI制御部と、

前記第2のLSI側とデータの送受信を行い、通信回線を介して前記他系の物理マシン側とデータの送受信を行う通信制御部と

を更に具備し、

前記第2のLSIは、

前記自系の物理マシン及び前記他系の物理マシンの電源状態を監視する電源監視部と、

前記自系の物理マシン及び前記他系の物理マシンの電源状態を制御する電源制御部と、

SMBusを介して、前記自系の物理マシン内部のBMC(Baseboard Management Controller)と接続し、前記BMCから監視結果を取得するSMBus制御部と、

前記第1のLSI側とデータの送受信を行い、通信回線を介して前記他系の物理マシン側とデータの送受信を行い、前記自系の物理マシンから電源の給電を受けられない場合、前記他系の物理マシンから電源の給電を受ける通信制御部と

を更に具備する

多重化システム。

【請求項5】

請求項1に記載の多重化システムであって、

前記データ通信カードは、

前記自系の物理マシンの基板上にチップとして搭載され、前記自系の物理マシン内部の前記基板上にそれぞれ配置されているCPUとI/Oチップとの間に設けられている

多重化システム。

【請求項6】

複数の物理マシンの各々に搭載されたデータ通信カードであって、

通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行う手段と、

メイン電源の給電により駆動し、前記自系の物理マシン及び前記他系の物理マシンの内部状態を監視し、状態異常を検出する第1のLSIと

を具備し、

前記第1のLSIは、

10

20

30

40

50

前記メイン電源がONの状態、監視及び検出対象の物理マシンのメモリに対してデータの読み出し要求を定期的に発行し、前記読み出し要求に対する応答を確認することで、前記監視及び検出対象の物理マシンの状態異常を検出する

データ通信カード。

【請求項7】

請求項6に記載のデータ通信カードであって、

前記データ通信カードは、

搭載された物理マシンの基板上にチップとして搭載され、前記基板上にそれぞれ配置されているCPUとI/Oチップとの間に設けられており、

前記第1のLSIは、

前記メイン電源がONの状態、前記監視及び検出対象の物理マシンの前記I/Oチップに対して設定情報の読み出し要求を定期的に発行し、前記読み出し要求に対する応答を確認することで、前記監視及び検出対象の物理マシンの状態異常を検出する

データ通信カード。

【請求項8】

請求項6または7に記載のデータ通信カードであって、

スタンバイ電源の給電により駆動し、前記メイン電源がOFFの状態の時でも、前記自系の物理マシン及び前記他系の物理マシンの電源状態を監視し、電源の異常を検出する第2のLSIと

を更に具備する

データ通信カード。

【請求項9】

請求項8に記載のデータ通信カードであって、

前記第1のLSIは、

前記データ通信カード内部を制御するプロセッサと、

前記自系の物理マシン上で動作するソフトウェア(SW)の状態を取得するSW状態取得部と、

PCI Expressバスを介して、前記自系の物理マシン内部のハードウェア及びI/Oチップと接続し、前記ハードウェア及び前記I/Oチップの状態を監視するPCI制御部と、

前記第2のLSI側とデータの送受信を行い、通信回線を介して前記他系の物理マシン側とデータの送受信を行う通信制御部と

を更に具備し、

前記第2のLSIは、

前記自系の物理マシン及び前記他系の物理マシンの電源状態を監視する電源監視部と、

前記自系の物理マシン及び前記他系の物理マシンの電源状態を制御する電源制御部と、

SMBusを介して、前記自系の物理マシン内部のBMC(Baseboard Management Controller)と接続し、前記BMCから監視結果を取得するSMBus制御部と、

前記第1のLSI側とデータの送受信を行い、通信回線を介して前記他系の物理マシン側とデータの送受信を行い、前記自系の物理マシンから電源の給電を受けられない場合、前記他系の物理マシンから電源の給電を受ける通信制御部と

を更に具備する

データ通信カード。

【請求項10】

複数の物理マシンの各々に搭載されたデータ通信カードにより実施される状態異常検出方法であって、

前記データ通信カードは、

通信手段と、

メイン電源の給電により駆動する第1のLSIと

10

20

30

40

50

を具備しており、
 前記状態異常検出方法は、
 前記通信手段が、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行うステップと、
 前記第1のLSIが、前記自系の物理マシン及び前記他系の物理マシンの内部状態を監視し、状態異常を検出するステップと
 を含み、

前記状態異常を検出するステップでは、
 前記第1のLSIが、前記メイン電源がONの状態、監視及び検出対象の物理マシンのメモリに対してデータの読み出し要求を定期的に発行し、前記読み出し要求に対する応答を確認することで、前記監視及び検出対象の物理マシンの状態異常を検出する
 状態異常検出方法。

10

【請求項11】

複数の物理マシンの各々に搭載されたデータ通信カードにより実行されるプログラムであって、

前記データ通信カードは、
 通信手段と、
 メイン電源の給電により駆動する第1のLSIと
 を具備しており、

20

前記プログラムは、
 前記通信手段が、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行うステップと、
 前記第1のLSIが、前記自系の物理マシン及び前記他系の物理マシンの内部状態を監視し、状態異常を検出するステップと
 を含み

前記状態異常を検出するステップでは、
 前記第1のLSIが、前記メイン電源がONの状態、監視及び検出対象の物理マシンのメモリに対してデータの読み出し要求を定期的に発行し、前記読み出し要求に対する応答を確認することで、前記監視及び検出対象の物理マシンの状態異常を検出すること
 をデータ通信カードに実行させるためのプログラム。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、多重化システムに関し、特にフォールトトレラントシステムやクラスターシステムを実現する多重化システム、データ通信カード、状態異常検出方法、及びプログラムに関する。

【背景技術】

【0002】

現在、物理マシン上で複数のOS(Operating System)を動作させることを可能とした仮想化技術が広く用いられている。仮想化技術によって、利用率の低い複数の仮想マシンを、1つの物理マシン上にまとめて、物理マシン1台あたりの利用効率を高めるとともに、物理マシンの台数を減らして消費電力を抑えることが可能となる。

40

【0003】

仮想化技術には、物理マシン上で稼働するホストOS上に仮想マシンを動作させる層を設け、その層上でゲストOSを動作させる方式や、ホストOSを介さず、ハードウェア(HW)上に仮想マシンを動作させるハイパーバイザ(hypervisor)を設け、そのハイパーバイザ上でゲストOSを動作させる方式等がある。

【0004】

また、現在、耐障害性の高いコンピュータシステムとして、フォールトトレラントシステム(Fault Tolerant System)が広く用いられている。

50

【 0 0 0 5 】

例えば、専用のハードウェア（HW）を介してロックステップ方式の動作を行い、多重化（多くは二重化）された主要なハードウェア（HW）を障害発生時に遅滞なく切り替えて動作し続けるハードウェア（HW）方式のフォールトトレラントシステムが、従来から知られている。

【 0 0 0 6 】

また、仮想マシンが動作している物理マシンでハードウェア（HW）上の故障等による障害が発生した場合に、その仮想マシンが行っている処理を、他の物理マシン上で待機している仮想マシンが継続して実行するソフトウェア（SW）方式のフォールトトレラントシステムが、近年研究されている。

10

【 0 0 0 7 】

ハードウェア（HW）方式のフォールトトレラントシステムでは、高価な専用のハードウェア（HW）を1システム毎に多重化する必要があるため、システムコストがかさむことになる。これに対して、ソフトウェア（SW）方式のフォールトトレラントシステムでは、汎用のハードウェア（HW）上で動作する複数の仮想マシンを用いて擬似的に多重化するため、専用のハードウェア（HW）が不要となり、システム毎にハードウェア（HW）を多重化する必要もないため、安価にシステムを構築・維持でき、システムコストを抑えることができる。

【 0 0 0 8 】

ソフトウェア（SW）方式のフォールトトレラントシステムでの処理の主体の切り替え方法の例として、ハードウェア（HW）方式のフォールトトレラントシステムで行われているロックステップ方式の動作や、メモリコピー方式の動作を、ソフトウェア（SW）により行い、障害発生時には瞬時に、処理の主体となる仮想マシンを切り替えるという方法がある。

20

【 0 0 0 9 】

しかし、一般的なIAサーバでソフトウェア（SW）方式のフォールトトレラントシステムやクラスターシステムを構築した場合、装置の故障検出を、一般的なハードウェア（HW）のエラー検出機構や専用のエラー検出ソフトウェア（SW）に依存することになる。

【 0 0 1 0 】

一般的なハードウェア（HW）では、機能を動作させた場合のみ故障を検出することができる。しかし、待機系となり動作していない場合には、故障を検出することができない。また、検出した故障を、ソフトウェア（SW）を介在せずに他系装置に通知する仕組みを持っていない。

30

【 0 0 1 1 】

また、専用のエラー検出ソフトウェア（SW）では、ソフトウェア（SW）によって動作させることができる機能については定期的なヘルスチェックを行うことで、早期に故障を検出したり、検出した装置を他系装置に通知したりできる。しかし、ソフトウェア（SW）が稼働するまでは、その機能を使用することができない。また、一般的に、ソフトウェア（SW）はタイムアウトによって故障を検出するため、他系装置に通知するには一定の時間を要する。したがって、フェイルオーバー（failover）時間を短くすることが困難である。

40

【 0 0 1 2 】

関連する技術として、特許文献1（特許第4468426号公報）に高可用システム及び実行状態制御方法が開示されている。この関連技術では、第1の仮想計算機を管理する第1のハイパーバイザが備える収集部が、第1の仮想計算機について発生した、第1の仮想計算機に対する入力を伴うイベントに関する同期情報を収集する。また、第2の仮想計算機を管理する第2のハイパーバイザが備える制御部が、この同期情報に従って、第2の仮想計算機の入力に係る実行状態を、第1の仮想計算機の入力に係る実行状態と同一になるように制御する。これにより、独立した2台の計算機上でそれぞれ稼働する仮想計算機を

50

組み合わせて二重化を実現する。

【0013】

また、特許文献2（特開2009-080692号公報）に仮想計算機システム及び同システムにおけるサービス引き継ぎ制御方法が開示されている。この関連技術では、仮想マシンが動作しているサーバ計算機に障害が発生した場合、サーバ計算機の仮想マシンモニタは、障害発生時刻に最も近い時点でディスク装置に採取されたスナップショットに基づき、仮想マシンを仮想マシンとしてサーバ計算機上に再生成する。通信記録ユニットの状態再現部は、仮想マシンに対応付けられた通信履歴に基づき、スナップショットの採取時期から上記障害発生時刻までの期間における仮想マシンの状態を仮想マシンに再現させる。再起動部は、例えば仮想マシンの状態の再現に失敗した場合、仮想マシンをサーバ計算機上で再起動する。これにより、仮想マシンが動作している物理計算機に障害が発生した場合、別の物理計算機上で再生成または再起動される仮想マシンによりサービスを継続させる。

10

【0014】

また、特許文献3（特開2008-033483号公報）に計算機システム、計算機および計算機動作環境の移動方法が開示されている。この関連技術では、第1計算機の動作を中断する。次に、第1ディスク上のコピーイメージに含まれるファイルのリストを作成する。次に、第1計算機の実行コンテキストを、第2計算機にコピーする。次に、第2計算機において上記動作を再開させる。次に、上記リストを参照して、コピーイメージを第1ディスクから第2ディスクにコピーする。これにより、第1ディスクを使用する第1計算機の動作環境を、第2ディスクを使用する第2計算機に移動させる際に、業務の中断時間を短縮する。

20

【先行技術文献】

【特許文献】

【0015】

【特許文献1】特許第4468426号公報

【特許文献2】特開2009-080692号公報

【特許文献3】特開2008-033483号公報

【発明の概要】

【発明が解決しようとする課題】

30

【0016】

従来の多重化システムでは、汎用装置を用いて故障検出を行っているが、汎用装置の故障検出能力は低い。例えば、故障検出機能を動作させた時にしか故障を検出できないため、待機系等で装置や故障検出機能が停止している場合は検出できない。また、多重化用OSの介在なく、外部に故障を通知することができないため、多重化した各装置のOSが起動していなければならない。更に、多重化用OSの故障検出機能も、多重化用OSが動作を開始するまでは機能しない。一般的に、ソフトウェア（SW）による故障検出では、タイムアウトで故障検出するため、検出まで時間がかかる。

【0017】

本発明の目的は、一般的なIAサーバでソフトウェア（SW）方式のフォールトトレラントシステムやクラスターシステムを構築するために利用されるデータ通信カードに、自系装置のメモリやチップ等の状態監視機能と、他系装置への状態通知機能及び電源制御機能を追加した多重化システムを提供することである。

40

【課題を解決するための手段】

【0018】

本発明に係る多重化システムは、複数の物理マシンと、複数の物理マシンの各々に搭載され、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行うためのデータ通信カードとを含む。データ通信カードは、自律的に、自系の物理マシン及び他系の物理マシンの状態を監視し、状態異常を検出する。

【0019】

50

本発明に係るデータ通信カードは、複数の物理マシンの各々に搭載されたデータ通信カードであって、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行う手段と、自律的に、自系の物理マシン及び他系の物理マシンの状態を監視し、状態異常を検出する手段とを具備する。

【0020】

本発明に係る状態異常検出方法は、複数の物理マシンの各々に搭載されたデータ通信カードにより実施される状態異常検出方法であって、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行うことと、自律的に、自系の物理マシン及び他系の物理マシンの状態を監視し、状態異常を検出することを含む。

【0021】

本発明に係るプログラムは、複数の物理マシンの各々に搭載されたデータ通信カードにより実行されるプログラムであって、通信回線を介して自系の物理マシンと他系の物理マシンとを接続し、相互にデータの送受信を行うステップと、自律的に、自系の物理マシン及び他系の物理マシンの状態を監視し、状態異常を検出するステップとをデータ通信カードに実行させるためのプログラムである。なお、本発明に係るプログラムは、記憶装置や記憶媒体に格納することが可能である。

【発明の効果】

【0022】

データ通信カードが、物理マシンの状態に関わらず、定期的に物理マシン内部の状態を常に監視するため、早期のエラー検出が可能となる。また、多重化システムにおけるスタンバイ系（待機系、従系）装置等、普段は動作していない装置のエラー検出も可能となる。

【図面の簡単な説明】

【0023】

【図1】本発明に係る多重化システムの基本構成を示す概念図である。

【図2】本発明に係る物理マシンの構成例を示す概念図である。

【図3】本発明に係るデータ通信カードの構成例を示す概念図である。

【図4】本発明に係るデータ通信カード内部の詳細を示すブロック図である。

【図5】本発明に係る多重化システムのシステム構成1を示す概念図である。

【図6】本発明に係る多重化システムのシステム構成2を示す概念図である。

【発明を実施するための形態】

【0024】

<実施形態>

以下に、本発明の実施形態について添付図面を参照して説明する。

【0025】

[基本構成]

図1に示すように、本発明に係る多重化システムは、複数の物理マシン100(100-i, i=1~n:nは台数)を含む。

【0026】

ここでは、物理マシン100(100-i, i=1~n)の各々の例として、PC(パソコン)、アプライアンス(appliance)、シンクライアントサーバ、ワークステーション、メインフレーム、スーパーコンピュータ等の計算機を想定している。なお、物理マシン100(100-i, i=1~n)は、端末やサーバに限らず、中継機器でも良い。中継装置の例として、ネットワークスイッチ(network switch)、ルータ(router)、プロキシ(proxy)、ゲートウェイ(gateway)、ファイアウォール(firewall)、ロードバランサ(load balancer: 負荷分散装置)、帯域制御装置(packet shaper)、セキュリティ監視制御装置(SCADA: Supervisory Control And Data Acquisition)、ゲートキーパー(gatekeeper)、基地局(base station)、アクセスポイント(AP: Access Point)、或いは、

10

20

30

40

50

複数の通信ポートを有する計算機等が考えられる。

【0027】

図示しないが、上記のような計算機や中継機器等は、プログラムに基づいて駆動し所定の処理を実行するプロセッサと、当該プログラムや各種データを記憶するメモリと、ネットワークインターフェースによって実現される。

【0028】

上記のプロセッサの例として、CPU (Central Processing Unit)、ネットワークプロセッサ (NP: Network Processor)、マイクロプロセッサ (microprocessor)、マイクロコントローラ、或いは、専用の機能を有する半導体集積回路 (IC: Integrated Circuit) 等が

10

【0029】

上記のメモリの例として、RAM (Random Access Memory)、ROM (Read Only Memory)、EEPROM (Electrically Erasable and Programmable Read Only Memory) やフラッシュメモリ等の半導体記憶装置等を想定している。HDD (Hard Disk Drive) やSSD (Solid State Drive) 等の補助記憶装置、又は、DVD (Digital Versatile Disk) 等のリムーバブルディスクや、SDメモリカード (Secure Digital memory card) 等の記憶媒体 (メディア) 等が考えられる。また、バッファ (buffer) やレジスタ (register) でも良い。或いは、DAS (Direct Attached Storage)、FC-SAN (Fibre Channel - Storage Area Network)、NAS (Network Attached Storage)、IP-SAN (IP - Storage Area Network) 等を用いたストレージ装置でも良い。

20

【0030】

なお、上記のプロセッサ及び上記のメモリは、一体化していても良い。例えば、近年では、マイコン等の1チップ化が進んでいる。したがって、電子機器等に搭載される1チップマイコンが、上記のプロセッサ及び上記のメモリを備えている事例が考えられる。

【0031】

上記のネットワークインターフェースの例として、ネットワーク通信に対応した基板 (マザーボード、I/Oボード) やチップ等の半導体集積回路、NIC (Network Interface Card) 等のネットワークアダプタや同様の拡張カード、アンテナ等の通信装置、通信ポート等が考えられる。

30

【0032】

また、ネットワークの例として、インターネット、LAN (Local Area Network)、無線LAN (Wireless LAN)、WAN (Wide Area Network)、バックボーン (Backbone)、ケーブルテレビ (CATV) 回線、固定電話網、携帯電話網、WiMAX (IEEE 802.16a)、3G (3rd Generation)、専用線 (lease line)、IrDA (Infrared Data Association)、Bluetooth (登録商標)、シリアル通信回線、データバス等が考えられる。

40

【0033】

但し、実際には、これらの例に限定されない。

【0034】

本発明では、物理マシン100 (100-i, i=1~n) の各々には、データ通信カード10が搭載されている。搭載の方法として、挿入、内蔵、組込、接続等が考えられる。但し、実際には、これらの例に限定されない。

【0035】

データ通信カード10は、物理マシン100 (100-i, i=1~n) に搭載可能な

50

高機能な拡張カードである。データ通信カード10は、上記のネットワークインターフェースでも良い。データ通信カード10は、LSI (Large Scale Integration) を搭載している。データ通信カード10は、他のデータ通信カードと通信回線を介して接続し、互いにデータの送受信を行う。また、データ通信カード10は、障害検出機能を持ち、自身が搭載された物理マシン内部の障害や、接続相手のデータ通信カードが搭載された他の物理マシン内部の障害を検出する。他の物理マシン内部の障害については、接続相手のデータ通信カードからの通知により認識しても良い。

【0036】

なお、データ通信カード10の形状は、カード型に限らない。例えば、データ通信カード10は、物理マシン100 (100 - i, i = 1 ~ n) の基板と一体化していても良い。但し、実際には、これらの例に限定されない。

10

【0037】

物理マシン100 (100 - i, i = 1 ~ n) 上に、データ通信カード10を使用できる環境が最低限整っていて、対応したソフトウェアが導入されている / 導入可能である場合、物理マシン100 (100 - i, i = 1 ~ n) にデータ通信カード10を挿すだけで、物理マシン100 (100 - i, i = 1 ~ n) は、FTサーバ (Fault Tolerant Server) やクラスターサーバになる。ソフトウェア (SW) 方式のフォールトトレラントシステムやクラスターシステムを実現する場合、データ通信カード10自体が、仮想マシン用の設定情報やイメージファイル等を保持し、物理マシン100 (100 - i, i = 1 ~ n) に提供するようにしていても良い。このとき、データ通信カード10は、「FTカード」や「クラスターカード」と呼ぶこともできる。

20

【0038】

[物理マシンの内部構成]

図2を参照して、物理マシン100 (100 - i, i = 1 ~ n) の各々の内部構成の詳細について説明する。

【0039】

物理マシン100 (100 - i, i = 1 ~ n) の各々は、データ通信カード10と、ハードウェア (HW: hardware) 20と、ソフトウェア (SW: software) 30と、ドライバ40と、I/O (Input/Output) チップ50と、BMC (Baseboard Management Controller) 60を備える。

30

【0040】

ここでは、2台の物理マシン (物理マシン100 - 1、物理マシン100 - 2) を例に説明する。例えば、物理マシン100 - 1は、データ通信カード10 - 1と、ハードウェア (HW) 20 - 1と、ソフトウェア (SW) 30 - 1と、ドライバ40 - 1と、BMC 60 - 1と、I/Oチップ50 - 1を備える。また、物理マシン100 - 2は、データ通信カード10 - 2と、ハードウェア (HW) 20 - 2と、ソフトウェア (SW) 30 - 2と、ドライバ40 - 2と、BMC 60 - 2と、I/Oチップ50 - 2を備える。

【0041】

データ通信カード10については、上記の通りである。

【0042】

ハードウェア (HW) 20は、物理マシン100 (100 - i, i = 1 ~ n) 内部のハードウェア (HW) である。通常、ハードウェア (HW) 20は、DC電源の給電 (電力供給) を受けて駆動する。ハードウェア (HW) 20の例として、プロセッサ、メモリ、補助記憶装置、ネットワークインターフェース、PCI (Peripheral Components Interconnect bus) スロット、及び電源装置、又はこれらの組み合わせ等が考えられる。なお、ハードウェア (HW) 20は、同一の物理マシン内部で多重化されていても良い。

40

【0043】

ソフトウェア (SW) 30は、ハードウェア (HW) 20を利用し、物理マシン100 (100 - i, i = 1 ~ n) 上で動作するソフトウェア (SW) である。ソフトウェア (

50

SW) 30の例として、OSやアプリケーションソフトウェア、ミドルウェア(middle ware)等が考えられる。なお、ソフトウェア(SW) 30は、物理マシン上に構築された仮想マシン(VM: Virtual Machine)でも良い。

【0044】

ドライバ40は、物理マシン100(100-i, i=1~n)の内部に装着された装置や、外部に接続した機器を制御・操作するためのソフトウェア(SW)/デバイスドライバである。ドライバ40は、OSが上記のような機器を制御するための橋渡しを行う。なお、ドライバ40は、OSに組み込まれ、OSの機能の一部として振舞うようにしても良い。すなわち、ドライバ40は、ソフトウェア(SW) 30の一部でも良い。ドライバ40は、ソフトウェア(SW) 30がI/Oチップ50に接続された外部のデバイスを利用する際、OSが提供する共通化されたAPI(アプリケーション・プログラミング・インターフェース)によってデバイスの機能を利用できるようにして、抽象化されたAPIとデバイスとの間の対応を受け持つ。

10

【0045】

I/Oチップ50は、物理マシン100(100-i, i=1~n)に搭載された接続口(物理ポート)であり、物理マシン本体と各種周辺機器に接続して、それらの機器とデータをやり取りするための入出力インターフェースである。入出力インターフェースの主な規格として、キーボードやマウスを接続するPS/2、モデムやプリンタなどと双方向で通信を行うシリアルインターフェースのRS232C、ハードディスクドライブ(HDD)等と双方向接続するパラレルインターフェースのSCSI、主に内蔵型HDD等と双方向接続するパラレルインターフェースのIDE、本体と周辺機器全般とを双方向接続するシリアルインターフェースであるUSB、次世代の高速なSCSI規格であるIEEE1394といった規格が知られている。なお、I/Oチップ50は、スーパーI/O(Super Input/Output)チップや、I/Oコントローラ・ハブ(ICH: I/O Controller Hub)でも良い。

20

【0046】

BMC60は、物理マシン100(100-i, i=1~n)内部に設けられたコントローラである。BMC60は、ハードウェア(HW) 20の状態を常時監視し、ハードウェア(HW)エラーの発生をOS等に通知する。具体的には、BMC60は、電源ユニットからの供給電圧や冷却ファンの回転数、プロセッサを含む各種パーツの温度、SCSIターミネータの電源電圧等を常時監視している。例え本体の電源がOFFでも、コンセントからの電源コードが電源ユニットに接続されている限り、BMC60には電力が供給されるため、BMC60は、本体の電源がOFF状態でも、ハードウェア(HW) 20の状態監視を継続する。

30

【0047】

[データ通信カードの詳細]

図3を参照して、データ通信カード10の詳細について説明する。

【0048】

ここでは、データ通信カード10に搭載されたLSIの例として、FPGA(Field Programmable Gate Array)と、CPLD(Complex Programmable Logic Device)を使用して説明する。なお、FPGAやCPLDは一例に過ぎない。実際には、他のLSIでも良い。

40

【0049】

データ通信カード10は、FPGA11と、CPLD12を備える。

【0050】

例えば、データ通信カード10-1は、FPGA11-1と、CPLD12-1を備える。また、データ通信カード10-2は、FPGA11-2と、CPLD12-2を備える。

【0051】

FPGA11は、第1のLSIである。FPGA11は、ハードウェア(HW) 20、

50

ソフトウェア (S W) 3 0、及び I / O チップ 5 0 の状態監視を行う。 F P G A 1 1 には、メイン電源 (データ通信カード 1 0 が受電した D C 電源) が給電されている。

【 0 0 5 2 】

C P L D 1 2 は、第 2 の L S I である。 C P L D 1 2 は、 B M C 6 0 及び電源装置の状態監視を行う。 C P L D 1 2 には、スタンバイ電源 (データ通信カード 1 0 が受電した A C 電源から作成される D C 電源) が給電されている。なお、スタンバイ電源とは、電源管理を行うため、常時一定出力を供給するための出力である。スタンバイ電源を出力する回路は、メイン電源を O F F にしても動作している。

【 0 0 5 3 】

なお、 F P G A 1 1 と C P L D 1 2 は、相互に通信可能である。

10

【 0 0 5 4 】

[F P G A 及び C P L D の詳細]

図 4 を参照して、 F P G A 1 1 及び C P L D 1 2 の詳細について説明する。

【 0 0 5 5 】

F P G A 1 1 は、プロセッサ 1 1 1 と、 S W 状態取得部 1 1 2 と、 P C I 制御部 1 1 3 と、通信制御部 1 1 4 を備える。

【 0 0 5 6 】

プロセッサ 1 1 1 は、データ通信カード 1 0 内部の各部の制御やデータの計算・加工 (演算処理) を行う。例えば、プロセッサ 1 1 1 は、データ通信カード 1 0 内部の R A M 等に記憶されたプログラムに基づいて駆動し、所定の処理を実行する。プロセッサ 1 1 1 は、物理マシン 1 0 0 (1 0 0 - i , i = 1 ~ n) 内部のハードウェア (H W) 2 0 に含まれる C P U 等のプロセッサではなく、データ通信カード 1 0 内部の F P G A 1 1 上に設けられている。また、プロセッサ 1 1 1 は、 S W 状態取得部 1 1 2、 P C I 制御部 1 1 3、及び通信制御部 1 1 4 の動作を変更することが可能である。

20

【 0 0 5 7 】

S W 状態取得部 1 1 2 は、ソフトウェア (S W) 3 0 の状態を取得する。 S W 状態取得部 1 1 2 は、ドライバ 4 0 から直接通知を受け取ることで、ソフトウェア (S W) 3 0 の状態を取得しても良いし、 P C I 制御部 1 1 3 を介して、ドライバ 4 0 がメモリ上に設定したソフトウェア (S W) 3 0 の状態を取得しても良い。

【 0 0 5 8 】

P C I 制御部 1 1 3 は、 P C I E x p r e s s バスを介して、物理マシン 1 0 0 (1 0 0 - i , i = 1 ~ n) 内部のハードウェア (H W) 2 0 及び I / O チップ 5 0 と接続し、ハードウェア (H W) 2 0 及び I / O チップ 5 0 の状態を監視する。

30

【 0 0 5 9 】

通信制御部 1 1 4 は、 F P G A 1 1 と C P L D 1 2 とを接続している。したがって、通信制御部 1 1 4 は、 F P G A 1 1 によるエラー検出結果を、 C P L D 1 2 側に通知することが可能である。また、通信制御部 1 1 4 は、他のデータ通信カードに搭載された F P G A と、少なくとも 1 本のケーブルを介してリモート (r e m o t e : 遠隔) 接続する。ここでは、 2 本のケーブルを介して接続することとする。ケーブルの本数が複数の場合、物理的 / 論理的に 1 本に束ねることも可能である。複数のケーブルを 1 本に束ねることで、束ねた本数に応じて通信速度を倍増することができる。また、複数のケーブルを 1 本に束ねることで、いずれかのケーブルが故障した場合でも、残りのケーブルを使用して通信を継続することが可能である。通信制御部 1 1 4 は、この 2 本のケーブルを介して、 F P G A 1 1 によるエラー検出結果を、他のデータ通信カードに通知する。

40

【 0 0 6 0 】

C P L D 1 2 は、電源監視部 1 2 1 と、電源制御部 1 2 2 と、 S M B u s (S y s t e m M a n a g e m e n t B u s) 制御部 1 2 3 と、通信制御部 1 2 4 を備える。

【 0 0 6 1 】

電源監視部 1 2 1 は、自機 (L o c a l : 自系) の物理マシン及びリモート接続先 (R e m o t e : 他系) の物理マシンの電源状態を監視する。このとき、電源監視部 1 2 1 は

50

、自機の物理マシンの電源状態を直接監視しても良いし、SMBus制御部123を介して、BMC60による監視結果を取得しても良い。なお、電源監視部121は、メイン電源の有無/変化により、自機の物理マシンの電源状態を監視しても良い。また、電源監視部121は、通信制御部124を介して、リモート接続先の物理マシンの電源状態を直接監視しても良いし、リモート接続先の物理マシンのSMBus制御部を介して、リモート接続先の物理マシンのBMCによる監視結果を取得しても良い。

【0062】

電源制御部122は、自機の物理マシン及びリモート接続先の物理マシンの電源状態を制御する。このとき、電源制御部122は、自機の物理マシンの電源状態を直接制御しても良いし、SMBus制御部123及びBMC60を介して制御しても良い。また、電源制御部122は、通信制御部124を介して、リモート接続先の物理マシンの電源状態を直接制御しても良いし、リモート接続先の物理マシンのSMBus制御部及びBMCを介して制御しても良い。

10

【0063】

SMBus制御部123は、SMBusを介して、物理マシン100(100-i, i=1~n)内部のBMC60と接続し、BMC60から監視結果を取得する。また、SMBus制御部123は、BMC60に対して制御指示を行うことも可能である。

【0064】

通信制御部124は、FPGA11とCPLD12とを接続している。したがって、通信制御部124は、CPLD12によるエラー検出結果を、FPGA11側に通知することが可能である。なお、FPGA11の通信制御部114とCPLD12の通信制御部124とは、同一の装置/回路でも良い。また、通信制御部124は、他のデータ通信カードに搭載されたCPLDと、少なくとも1本のケーブルを介してリモート接続する。ここでは、1本のケーブルを介して接続することとする。ケーブルの本数が複数の場合、物理的/論理的に1本に束ねることが可能である。通信制御部124は、この1本のケーブルを介して、CPLD12によるエラー検出結果を、他のデータ通信カードに通知する。また、通信制御部124は、この1本のケーブルを介して、リモート接続先の物理マシンから電源の給電を受ける(受電する)ようにしても良い。例えば、通信制御部124は、この1本のケーブルを構成する回線の一部を、リモート接続先の物理マシンからの電源の給電用に利用する。この場合、CPLD12は、自機の物理マシンから電源の給電を受けることができなくなっても、リモート接続先の物理マシンから電源の給電を受けることで、動作を継続することができる。

20

30

【0065】

[ハードウェア(HW)状態監視]

以下に、ハードウェア(HW)状態監視の動作の詳細について説明する。

【0066】

ここでは、ハードウェア(HW)状態監視の一例として、メモリ状態監視の動作について説明する。なお、メモリ状態監視が行われるのは、メイン電源がON状態(メイン電源とスタンバイ電源の両方がONの状態)の時のみである。

【0067】

PCI制御部113は、PCI Expressバスを介して、メモリの読み出し要求(リードリクエスト)を発行し、CPU等を介してメモリのデータの読み出しを行い、CPU等からコンプリーション(completion:完了)の応答が正常に返ってくるかをチェックすることで、メモリが正常に動作しているか監視する。

40

【0068】

例えば、PCI制御部113は、一定間隔で定期的に、全メモリアドレスに対して順番に読み出し要求(リードリクエスト)を発行する。

【0069】

なお、PCI制御部113の動作は、プロセッサ111により変更することが可能である。

50

【 0 0 7 0 】

これにより、P C I 制御部 1 1 3 は、メモリ故障を検出することが可能となる。

【 0 0 7 1 】

従来は、実際にメモリの読み出しが必要となり、C P U 等からメモリの読み出し要求（リードリクエスト）が発行されるまで、エラー検出が不可能であった。

【 0 0 7 2 】

本発明では、データ通信カード 1 0 が、メモリの読み出しの有無に関わらず、定期的にメモリの読み出し要求（リードリクエスト）を発行し、メモリの状態を常に監視するため、早期のエラー検出が可能となる。

【 0 0 7 3 】

また、多重化システムにおけるスタンバイ系（待機系、従系）装置等、普段は動作していない装置のエラー検出も可能となる。

【 0 0 7 4 】

[ソフトウェア（S W）状態監視]

以下に、ソフトウェア（S W）状態監視の動作の詳細について説明する。

【 0 0 7 5 】

S W 状態取得部 1 1 2 は、ソフトウェア（S W）3 0 の状態を取得することで、ソフトウェア（S W）3 0 が正常に動作しているか監視する。

【 0 0 7 6 】

S W 状態取得部 1 1 2 は、ドライバ 4 0 から直接通知を受け取ることで、ソフトウェア（S W）3 0 の状態を取得しても良いし、P C I 制御部 1 1 3 を介して、ドライバ 4 0 がメモリ上に設定したソフトウェア（S W）3 0 の状態を取得しても良い。

【 0 0 7 7 】

これにより、S W 状態取得部 1 1 2 は、ソフトウェア（S W）3 0 の異常を検出することが可能となる。

【 0 0 7 8 】

通常、故障監視をしているソフトウェア（S W）自体を用いて、当該ソフトウェア（S W）の異常を検出することは困難である。

【 0 0 7 9 】

本発明では、データ通信カード 1 0 が、ソフトウェア（S W）3 0 の状態を監視するため、故障監視をしているソフトウェア（S W）自体のエラー検出が可能となる。

【 0 0 8 0 】

[I / O チップ状態監視機能]

以下に、I / O チップ状態監視機能の動作の詳細について説明する。

【 0 0 8 1 】

I / O チップ状態監視が行われるのは、メイン電源が O N 状態（メイン電源とスタンバイ電源の両方が O N の状態）の時のみである。

【 0 0 8 2 】

P C I 制御部 1 1 3 は、P C I E x p r e s s バスを介して、I / O チップ 5 0 の設定情報（コンフィグ）の読み出し要求（リードリクエスト）を発行し、直接に / C P U 等を介して I / O チップ 5 0 の設定情報の読み出しを行い、I / O チップ 5 0 / C P U 等からコンプリーション（c o m p l e t i o n : 完了）の応答が正常に返ってくるかをチェックすることで、I / O チップ 5 0 が正常に動作しているか監視する。

【 0 0 8 3 】

例えば、P C I 制御部 1 1 3 は、一定間隔で定期的に、全 I / O チップ 5 0 に対して順番に読み出し要求（リードリクエスト）を発行する。

【 0 0 8 4 】

なお、P C I 制御部 1 1 3 の動作は、プロセッサ 1 1 1 により変更することが可能である。

【 0 0 8 5 】

10

20

30

40

50

これにより、P C I 制御部 1 1 3 は、I / O チップ 5 0 の故障を検出することが可能となる。

【 0 0 8 6 】

従来は、実際に I / O チップ 5 0 の読み出しが必要となり、C P U 等から I / O チップ 5 0 の読み出し要求（リードリクエスト）が発行されるまで、エラー検出が不可能であった。

【 0 0 8 7 】

本発明では、データ通信カード 1 0 が、I / O チップ 5 0 の読み出しの有無に関わらず、定期的に I / O チップ 5 0 の読み出し要求（リードリクエスト）を発行し、I / O チップ 5 0 の状態を常に監視するため、早期のエラー検出が可能となる。

10

【 0 0 8 8 】

また、多重化システムにおけるスタンバイ系（待機系、従系）装置等、普段は動作していない装置のエラー検出も可能となる。

【 0 0 8 9 】

[B M C 状態監視]

以下に、B M C 状態監視の動作の詳細について説明する。

【 0 0 9 0 】

なお、B M C 状態監視は、メイン電源が O N 状態（メイン電源とスタンバイ電源の両方が O N の状態）/ メイン電源が O F F 状態（メイン電源が O F F の状態で、スタンバイ電源のみ O N の状態）のいずれの状態であっても行われる。

20

【 0 0 9 1 】

S M B u s 制御部 1 2 3 は、S M B u s を介して、B M C 6 0 に S M B u s 読み出し要求（リードリクエスト）を発行し、B M C 6 0 が持っているレジスタの値の読み出しを行い、B M C 6 0 からコンプリーション（c o m p l e t i o n : 完了）の応答が正常に返ってくるかをチェックすることで、B M C 6 0 が正常に動作しているか監視する。

【 0 0 9 2 】

例えば、S M B u s 制御部 1 2 3 は、一定間隔で定期的に、B M C 6 0 に S M B u s 読み出し要求（リードリクエスト）を発行する。

【 0 0 9 3 】

B M C 6 0 は、S M B u s 制御部 1 2 3 に対して、S M B u s 書き込み要求（ライトリクエスト）を発行し、データ通信カード 1 0 内にあるレジスタに対してデータの書き込みを行うことも可能である。S M B u s 制御部 1 2 3 は、B M C 6 0 から一定間隔で S M B u s 書き込み要求（ライトリクエスト）が発行されるかどうかをチェックすることで、B M C 6 0 が正常に動作しているか監視しても良い。

30

【 0 0 9 4 】

なお、メイン電源が O N 状態であれば、プロセッサ 1 1 1 は、S M B u s 制御部 1 2 3 の動作を変更することが可能である。

【 0 0 9 5 】

これにより、S M B u s 制御部 1 2 3 は、B M C 6 0 の故障を検出することが可能となる。

40

【 0 0 9 6 】

本発明では、データ通信カード 1 0 が、定期的に B M C 6 0 の S M B u s 読み出し要求（リードリクエスト）を発行し、B M C 6 0 の状態を常に監視するため、早期のエラー検出が可能となる。

【 0 0 9 7 】

[自機の電源状態監視]

以下に、自機の電源状態監視の動作の詳細について説明する。

【 0 0 9 8 】

なお、自機の電源状態監視は、メイン電源が O N 状態（メイン電源とスタンバイ電源の両方が O N の状態）/ メイン電源が O F F 状態（メイン電源が O F F の状態で、スタンバ

50

イ電源のみONの状態)のいずれの状態であっても行われる。

【0099】

電源監視部121は、自機の物理マシンの電源状態を監視する。なお、電源監視部121は、スタンバイ電源の給電により駆動している。

【0100】

(1)メイン電源がON状態時の動作

電源監視部121は、自機の物理マシンのメイン電源がON状態であることを検出する。また、電源監視部121は、通信制御部124を介して、リモート接続先の物理マシンのメイン電源がON状態であることを検出する。

【0101】

電源制御部122は、リモート接続先の物理マシンのメイン電源がOFF状態であることを検出した場合、通信制御部124を介して、リモート接続先の物理マシンのメイン電源をON状態にする。

【0102】

なお、自機の物理マシンのメイン電源がON状態であることを検出した場合、電源制御部122は、通信制御部124を介して、リモート接続先の物理マシンのメイン電源をON状態にすることが可能である。リモート接続先の物理マシンのメイン電源をON状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとを同時に起動する場合に有用である。

【0103】

例えば、ソフトウェア(SW)方式のフォールトトレラントシステムを運用する場合、必ず2台の物理マシン(アクティブ系、スタンバイ系)のDC電源がON状態になっている必要がある。そのため、1台の装置の電源スイッチを押した際、連動してもう片方の装置もDC電源をON状態にすることに利点(メリット)がある。逆に、1台の装置を停止(Shutdown)してDC電源をOFF状態にする際、連動してもう片方の装置もDC電源をOFF状態にすることに利点がある。2台の物理マシンを連動させるかどうかは固定的ではなく、選択可能である。

【0104】

また、電源制御部122は、自機の物理マシンのメイン電源がON状態であることを検出した場合、リモート接続先の物理マシンのメイン電源をOFF状態にすることも可能である。リモート接続先の物理マシンのメイン電源をOFF状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとの間で、アクティブ系(実行系、主系)とスタンバイ系(待機系、従系)を切り替える場合に有用である。

【0105】

(2)メイン電源がOFF状態時の動作

電源監視部121は、自機の物理マシンのメイン電源がOFF状態であることを検出する。また、電源監視部121は、通信制御部124を介して、リモート接続先の物理マシンのメイン電源がOFF状態であることを検出する。

【0106】

電源制御部122は、リモート接続先の物理マシンのメイン電源がON状態であることを検出した場合、通信制御部124を介して、リモート接続先の物理マシンのメイン電源をOFF状態にする。

【0107】

なお、自機の物理マシンのメイン電源がOFF状態であることを検出した場合、電源制御部122は、通信制御部124を介して、リモート接続先の物理マシンのメイン電源をOFF状態にすることが可能である。リモート接続先の物理マシンのメイン電源をOFF状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとを同時に停止する場合に有用である。

【0108】

10

20

30

40

50

また、電源制御部 1 2 2 は、自機の物理マシンのメイン電源が OFF 状態であることを検出した場合、リモート接続先の物理マシンのメイン電源を ON 状態にすることも可能である。リモート接続先の物理マシンのメイン電源を ON 状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとの間で、アクティブ系（実行系、主系）とスタンバイ系（待機系、従系）を切り替える場合に有用である。

【 0 1 0 9 】

これにより、電源監視部 1 2 1 は、電源故障を検出することが可能となる。

【 0 1 1 0 】

通常、ソフトウェア（SW）制御により物理マシンの電源故障を検出することは困難である。 10

【 0 1 1 1 】

自機の物理マシンのソフトウェア（SW）により自機の物理マシンの電源故障を検出することが困難な理由は、自機の物理マシンの電源が故障すると、自機の物理マシンのソフトウェア（SW）が停止する可能性が高いためである。

【 0 1 1 2 】

本発明では、データ通信カード 1 0 が、自機の物理マシン及びリモート接続先の物理マシンの電源状態を監視するため、各物理マシンの電源故障の検出、及びいずれの物理マシンの電源故障であるかの判別・特定が可能となる。

【 0 1 1 3 】

[リモート接続先の電源状態監視]

以下に、リモート接続先の電源状態監視の動作の詳細について説明する。 20

【 0 1 1 4 】

なお、リモート接続先の電源状態監視は、メイン電源が ON 状態（メイン電源とスタンバイ電源の両方が ON の状態）/ メイン電源が OFF 状態（メイン電源が OFF の状態、スタンバイ電源のみ ON の状態）のいずれの状態であっても行われる。

【 0 1 1 5 】

電源監視部 1 2 1 は、リモート接続先の物理マシンの電源状態を監視する。なお、電源監視部 1 2 1 は、スタンバイ電源の給電により駆動している。 30

【 0 1 1 6 】

（ 1 ）メイン電源が ON 状態時の動作

電源監視部 1 2 1 は、通信制御部 1 2 4 を介して、リモート接続先の物理マシンのメイン電源が ON 状態であることを検出する。また、電源監視部 1 2 1 は、自機の物理マシンのメイン電源が ON 状態であることを検出する。

【 0 1 1 7 】

電源制御部 1 2 2 は、自機の物理マシンのメイン電源が OFF 状態であることを検出した場合、自機の物理マシンのメイン電源を ON 状態にする。

【 0 1 1 8 】

なお、リモート接続先の物理マシンのメイン電源が ON 状態であることを検出した場合、電源制御部 1 2 2 は、通信制御部 1 2 4 を介して、自機の物理マシンのメイン電源を ON 状態にすることが可能である。自機の物理マシンのメイン電源を ON 状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとを同時に起動する場合に有用である。 40

【 0 1 1 9 】

また、電源制御部 1 2 2 は、リモート接続先の物理マシンのメイン電源が ON 状態であることを検出した場合、自機の物理マシンのメイン電源を OFF 状態にすることも可能である。自機の物理マシンのメイン電源を OFF 状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとの間で、アクティブ系（実行系、主系）とスタンバイ系（待機系、従系）を切り替える場合に有用である。

【 0 1 2 0 】

(2) メイン電源がOFF状態時の動作

電源監視部121は、通信制御部124を介して、リモート接続先の物理マシンのメイン電源がOFF状態であることを検出する。また、電源監視部121は、自機の物理マシンのメイン電源がOFF状態であることを検出する。

【0121】

電源制御部122は、自機の物理マシンのメイン電源がON状態であることを検出した場合、自機の物理マシンのメイン電源をOFF状態にする。

【0122】

なお、リモート接続先の物理マシンのメイン電源がOFF状態であることを検出した場合、電源制御部122は、通信制御部124を介して、自機の物理マシンのメイン電源をOFF状態にすることが可能である。自機の物理マシンのメイン電源をOFF状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとを同時に停止する場合に有用である。

10

【0123】

また、電源制御部122は、リモート接続先の物理マシンのメイン電源がOFF状態であることを検出した場合、自機の物理マシンのメイン電源をON状態にすることも可能である。自機の物理マシンのメイン電源をON状態にするか否かについては、設定により変更可能である。自機の物理マシンとリモート接続先の物理マシンとの間で、アクティブ系（実行系、主系）とスタンバイ系（待機系、従系）を切り替える場合に有用である。

【0124】

これにより、電源監視部121は、電源故障を検出することが可能となる。

20

【0125】

通常、ソフトウェア（SW）制御により物理マシンの電源故障を検出することは困難である。

【0126】

自機の物理マシンのソフトウェア（SW）により自機の物理マシンの電源故障を検出することが困難な理由は、自機の物理マシンの電源が故障すると、自機の物理マシンのソフトウェア（SW）が停止する可能性が高いためである。

【0127】

自機の物理マシンのソフトウェア（SW）によりリモート接続先の物理マシンの電源故障を検出することが困難な理由は、リモート接続先の物理マシンの電源が壊れたのか、リモート接続先の物理マシンのソフトウェア（SW）の通信が止まったただけかを判断するのが困難なためである。

30

【0128】

本発明では、データ通信カード10が、自機の物理マシン及びリモート接続先の物理マシンの電源状態を監視するため、早期のエラー検出が可能となる。

【0129】

[自律制御によるメイン電源のOFF]

以下に、自律制御によりメイン電源をOFF状態にする動作の詳細について説明する。

【0130】

データ通信カード10は、上記の各状態監視の結果、いずれかの故障を検出した場合、ソフトウェア（SW）30の介在なく、自律的に、故障箇所に係る装置のメイン電源をOFF状態にする。このとき、データ通信カード10は、自機の物理マシン又はリモート接続先の物理マシン自体のメイン電源をOFF状態にしても良い。

40

【0131】

例えば、電源制御部122は、自機の物理マシン自体のメイン電源をOFF状態にすべきであれば、自機の物理マシン自体のメイン電源をOFF状態にする。

【0132】

また、電源制御部122は、リモート接続先の物理マシン自体のメイン電源をOFF状態にすべきであれば、通信制御部124を介して、リモート接続先の物理マシン自体のメ

50

イン電源をOFF状態にする。

【0133】

本発明では、データ通信カード10が、自律的にメイン電源をOFF状態にする制御を行うため、エラー検出後、ソフトウェア(SW)の介在なく、物理マシン自体のメイン電源をOFF状態にすることが可能となる。

【0134】

[自律制御によるメイン電源のON]

以下に、自律制御によりメイン電源をON状態にする動作の詳細について説明する。

【0135】

データ通信カード10は、リモート接続先の物理マシンにおける故障を検出した場合、自機の物理マシンのメイン電源がOFF状態であれば、自機の物理マシンのメイン電源をON状態にする。

10

【0136】

データ通信カード10は、通信制御部124を介して、リモート接続先の物理マシンにおける故障を検出する。

【0137】

電源制御部122は、リモート接続先の物理マシンにおける故障が発生した場合、自機の物理マシンのメイン電源をON状態にする。

【0138】

本発明では、データ通信カード10が、自律的にリモート接続先の物理マシンにおける故障を検出するため、リモート接続先の物理マシンのエラー検出後、ソフトウェア(SW)の介在なく、自機の物理マシンのメイン電源をON状態にすることが可能となる。

20

【0139】

これにより、アクティブ・スタンバイ方式の二重化システムにおいて、スタンバイ系(待機系、従系)装置のメイン電源をON状態にする必要がある時まで、メイン電源をOFF状態のまま待機させておくことが可能となる。したがって、スタンバイ系(待機系、従系)装置のメイン電源をON状態で待機させておく必要がなくなり、システム全体の消費電力を大幅に削減することが可能となる。

【0140】

[システム構成1(データ通信カード独立型)]

図5を参照して、データ通信カード10が、物理マシン100(100-i, i=1~n)の基板から独立して存在している「システム構成1」について説明する。

30

【0141】

ここでは、データ通信カード10は、物理マシン100(100-i, i=1~n)のカードスロットに挿された拡張カードである。なお、データ通信カード10の形状は、カード型に限らない。

【0142】

データ通信カード10は、ハードウェア(HW)20の1つであるプロセッサ(CPU等)を介して、I/Oチップ50と接続する。

【0143】

例えば、データ通信カード10は、PCI Expressバスを介して、物理マシン100(100-i, i=1~n)内部のハードウェア(HW)20の1つであるプロセッサ(CPU等)と接続する。このプロセッサ(CPU等)は、PCI Expressバスを介して、I/Oチップ50と接続する。

40

【0144】

また、データ通信カード10は、SMBusを介して、BMC60と接続する。

【0145】

[システム構成2(データ通信カード一体型)]

図6を参照して、データ通信カード10が、物理マシン100(100-i, i=1~n)の基板と一体化している「システム構成2」について説明する。

50

【0146】

ここでは、データ通信カード10は、物理マシン100(100-i, i=1~n)の基板に搭載されたチップである。この場合、物理マシン100(100-i, i=1~n)の基板自体が、データ通信カード10としての機能も持つことになる。すなわち、物理マシン100(100-i, i=1~n)の基板自体が、データ通信カード10に相当する。

【0147】

データ通信カード10は、ハードウェア(HW)20の1つであるプロセッサ(CPU等)とI/Oチップ50との間に存在し、このプロセッサ(CPU等)とI/Oチップ50との間の通信を監視する。

10

【0148】

例えば、物理マシン100(100-i, i=1~n)内部のハードウェア(HW)20の1つであるプロセッサ(CPU等)は、PCI Expressバスを介して、データ通信カード10と接続する。データ通信カード10は、PCI Expressバスを介して、I/Oチップ50と接続する。

【0149】

また、データ通信カード10は、SMBusを介して、BMC60と接続する。

【0150】

<本発明の特徴>

以上のように、本発明は、一般的なIAサーバでソフトウェア方式のフォールトトレラントシステムやクラスターシステムを構築するために利用されるデータ通信カードに、スタンバイ電源から動作可能な装置の状態監視機能、他系装置への状態通知機能及び電源制御機能を追加することで、早期の故障検出とフェイルオーバー(failover)及びコールドスタンバイ(cold standby)を実現する。

20

【0151】

上記の機能を追加したデータ通信カードを一般的なIAサーバに挿入するだけで、このデータ通信カードが自律的、定期的に、IAサーバの主要コンポーネントが正常に動作しているか否かチェックを行うため、早期の故障検出が可能となる。また、このデータ通信カードは、検出した故障を他系装置に通知し、他系装置を即座にフェイルオーバー処理に遷移させることができる。

30

【0152】

また、このデータ通信カードを使用してアクティブ・スタンバイ構成の多重化システムを構築した場合、スタンバイ側の装置は、メイン電源がOFF状態であっても他系装置の状態を監視可能であり、他系装置の故障を検出した場合、自律的にメイン電源がON状態にすることも可能である。

【0153】

このように、このデータ通信カードは、早期の故障検出とフェイルオーバーを実現し、第三者を介さない自律的なコールドスタンバイも実現する。

【0154】

<備考>

以上、本発明の実施形態を詳述してきたが、実際には、上記の実施形態に限られるものではなく、本発明の要旨を逸脱しない範囲の変更があっても本発明に含まれる。

40

【符号の説明】

【0155】

10... データ通信カード

11... FPGA(Field Programmable Gate Array)

111... プロセッサ

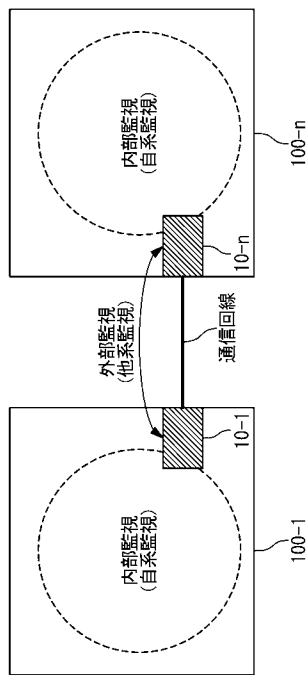
112... SW(software)状態取得部

113... PCI(Peripheral Components Interco

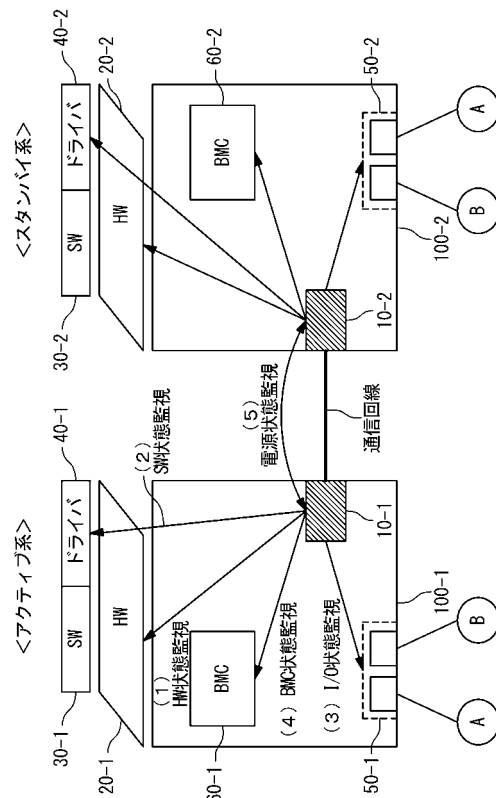
50

- connect bus) 制御部
- 114... 通信制御部
- 12... CPLD (Complex Programmable Logic Device)
- 121... 電源監視部
- 122... 電源制御部
- 123... SMBus (System Management Bus) 制御部
- 124... 通信制御部
- 20... ハードウェア (HW: hardware)
- 30... ソフトウェア (SW: software)
- 40... ドライバ
- 50... I/O (Input/Output) チップ
- 60... BMC (Baseboard Management Controller)
- r)
- 100(-i, i = 1 ~ n) ... 物理マシン

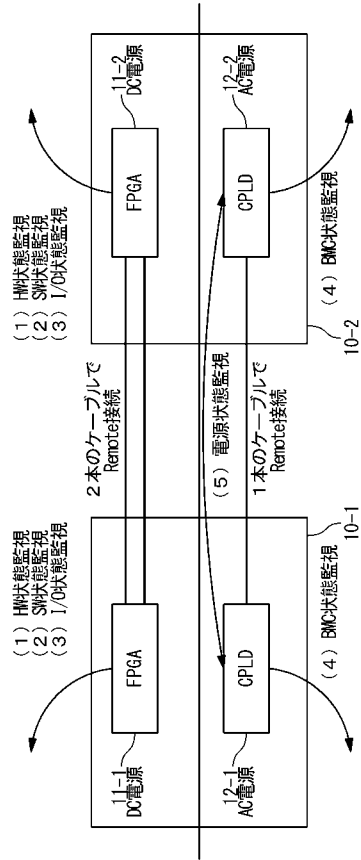
【図1】



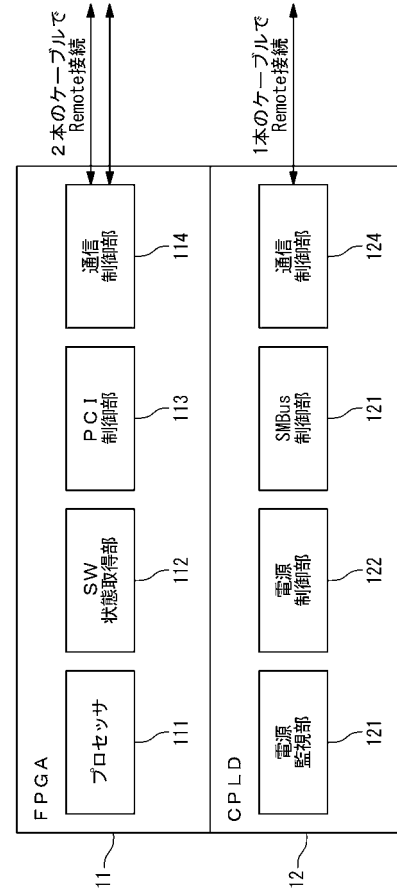
【図2】



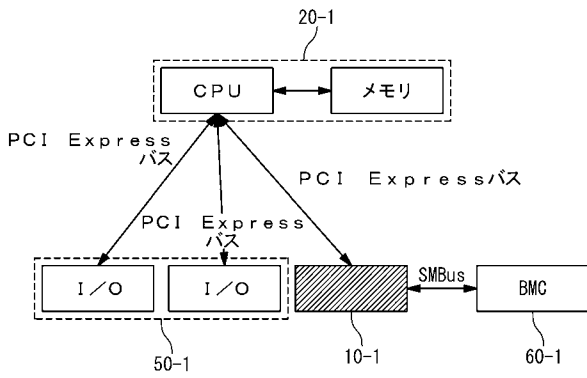
【 図 3 】



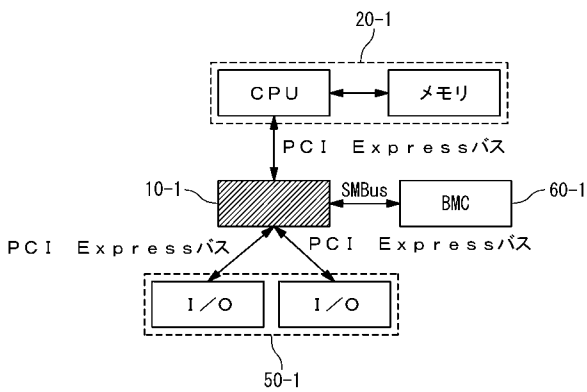
【 図 4 】



【 図 5 】



【 図 6 】



フロントページの続き

- (56)参考文献 特開2001-060160(JP,A)
特開2006-172050(JP,A)
特開2009-205334(JP,A)
特開2008-225567(JP,A)
特開平10-150479(JP,A)
特開2011-022957(JP,A)
特開2011-134314(JP,A)
特開平08-251814(JP,A)
特開2008-033483(JP,A)
特開2009-080692(JP,A)
特開2009-080695(JP,A)
特開2010-079789(JP,A)
特開2010-020505(JP,A)

(58)調査した分野(Int.Cl., DB名)

- G06F 11/16 - 11/20,
G06F 11/28 - 11/34,
G06F 1/26, 1/28