

(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2014-170394
(P2014-170394A)

(43) 公開日 平成26年9月18日(2014.9.18)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 11/20 (2006.01)	G06F 11/20 310E	5B034
G06F 3/06 (2006.01)	G06F 3/06 301A	
	G06F 3/06 306B	

審査請求 未請求 請求項の数 10 O L (全 21 頁)

(21) 出願番号 特願2013-42016 (P2013-42016)
(22) 出願日 平成25年3月4日 (2013.3.4)

(71) 出願人 000004237
日本電気株式会社
東京都港区芝五丁目7番1号
(74) 代理人 100124811
弁理士 馬場 資博
(74) 代理人 100088959
弁理士 境 廣巳
(72) 発明者 丹生谷 吉紀
東京都港区芝五丁目7番1号 日本電気株式会社内
Fターム(参考) 5B034 BB02 CC01 DD05

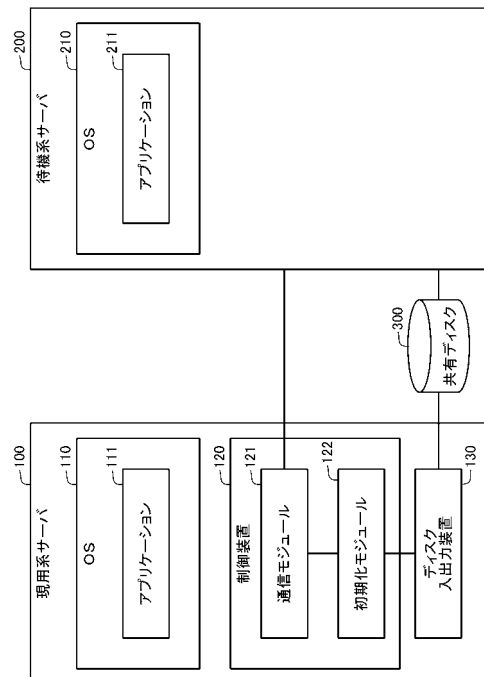
(54) 【発明の名称】 クラスタシステム

(57) 【要約】

【課題】共有データを保護しつつ、迅速にフェイルオーバーを行うこと。

【解決手段】本発明のクラスタシステムは、フェイルオーバー機能を有する現用系サーバ及び待機系サーバと、共有ディスクと、を備え、上記現用系サーバは、OSの影響を受けずに作動する制御装置と、共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えており、上記現用系サーバの制御装置は、待機系サーバと通信を行う通信モジュールと、現用系サーバの障害発生時にディスク入出力装置を初期化して通信モジュールを介して待機系サーバに通知する初期化モジュールと、を備えた。

【選択図】 図10



【特許請求の範囲】**【請求項 1】**

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備え、

前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えており、

前記現用系サーバの前記制御装置は、前記待機系サーバと通信を行う通信モジュールと、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールと、を備え、

前記待機系サーバは、前記現用系サーバの前記制御装置から前記 I/O フェンシングが完了した旨の通知を受けて、前記待機系サーバ上における前記アプリケーションにて前記現用系サーバ上における前記アプリケーションによる処理を引き継ぐフェイルオーバー処理機能を有する、

クラスタシステム。

【請求項 2】

請求項 1 に記載のクラスタシステムであって、

前記現用系サーバは、当該現用系サーバにて作動しているソフトウェアの障害発生を検出して前記制御装置の前記通知モジュールに通知する障害検出モジュールを備え、

前記制御装置の前記初期化モジュールは、前記障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、

クラスタシステム。

【請求項 3】

請求項 1 又は 2 に記載のクラスタシステムであって、

前記現用系サーバの前記制御装置は、当該現用系サーバに装備されたハードウェアの障害発生を検出して前記制御装置の前記初期化モジュールに通知するハードウェア障害検出モジュールを備え、

前記制御装置の前記初期化モジュールは、前記ハードウェア障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、

クラスタシステム。

【請求項 4】

請求項 1 乃至 3 のいずれかに記載のクラスタシステムであって、

前記待機系サーバは、前記現用系サーバにおける障害発生を検出して、前記現用系サーバの前記初期化モジュールに通知する待機系側障害検出モジュールを備え、

前記現用系サーバの前記初期化モジュールは、前記待機系サーバの前記待機系側障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、

クラスタシステム。

【請求項 5】

請求項 1 乃至 4 のいずれかに記載のクラスタシステムであって、

前記初期化モジュールは、前記ディスク入出力装置にアクセスするためのアクセス情報

10

20

30

40

50

と、前記待機系サーバに情報を通知するためのアドレス情報と、を記憶保持しており、前記アクセス情報に基づいて前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記アドレス情報に基づいて前記通信モジュールを介して前記待機系サーバに通知する、クラスタシステム。

【請求項6】

請求項5に記載のクラスタシステムであって、

前記現用系サーバは、当該現用系サーバの起動時に、前記アクセス情報及び前記アドレス情報が予め設定された設定ファイルを読み出し、当該設定ファイルに設定されている前記アクセス情報及び前記アドレス情報を前記制御装置の前記初期化モジュールに通知し、

10

前記初期化モジュールは、通知された前記アクセス情報及び前記アドレス情報を記憶保持する、

クラスタシステム。

【請求項7】

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムにおける前記現用系サーバであって、

前記現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えており、

20

前記制御装置は、前記待機系サーバと通信を行う通信モジュールと、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールと、を備えた、現用系サーバ。

【請求項8】

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムにおける前記現用系サーバが、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

30

前記制御装置に、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールによる処理、

を実行させるためのプログラム。

【請求項9】

40

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムによるフェイルオーバー方法であって、

前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記現用系サーバの前記制御装置は、前記現用系サーバの障害発生時に、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行

50

、当該 I / O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知し、

前記待機系サーバは、前記現用系サーバの前記制御装置から前記 I / O フェンシングが完了した旨の通知を受けて、前記待機系サーバ上における前記アプリケーションにて前記現用系サーバ上における前記アプリケーションによる処理を引き継ぐ、
フェイルオーバー方法。

【請求項 10】

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えた
クラスタシステムによるフェイルオーバー方法であって、

前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記現用系サーバの前記制御装置は、前記現用系サーバの障害発生時に、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I / O フェンシングを行い、当該 I / O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、

フェイルオーバー方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、クラスタシステムにかかり、特に、フェイルオーバー機能を有するクラスタシステムに関する。

【背景技術】

【0002】

コンピュータの障害発生による処理の停止を抑制するために、例えば、特許文献 1 に示すように、現用系及び待機系の両サーバと共有ディスクを備えたフェイルオーバー型クラスタシステムが知られている。このようなフェイルオーバー型のクラスタシステムでは、両サーバにて作動する業務アプリケーションのデータは、現用系サーバからも待機系サーバからもアクセス可能な共有ディスクに格納される。そして、現用系サーバに障害が発生した場合、業務アプリケーションの切り替え（以降「フェイルオーバー」と記す）を行うことで、クラスタソフトウェアは待機系サーバで業務アプリケーションを起動する。待機系サーバで起動した業務アプリケーションは、共有ディスクに格納されたデータを使用し、現用系サーバの業務アプリケーションが停止した時点から業務処理を再開する。

【0003】

上述したように、共有ディスク構成のフェイルオーバー型クラスタシステムでは、データを共有ディスクに格納するため、現用系と待機系の両サーバから同時に共有ディスクに書込を行うと、データが破壊するおそれがある。そのため、通常は現用系サーバからの書込のみを行うように排他制御を行う。以下に排他制御を行う例を記載する。

【0004】

（1：フェイルオーバー時の排他制御）

現用系サーバから待機系サーバにフェイルオーバーを実行する場合は、現用系サーバからの共有ディスクアクセスを確実に停止した後に、待機系サーバで業務アプリケーションを起動する必要がある。

【0005】

（2：スプリットブレイン状態の排他制御）

サーバ間のネットワークの障害によりクラスタがお互いの状態を認識できないスプリッ

10

20

30

40

50

トブレイン状態になった場合は、共有ディスクへのアクセスの有無により、業務アプリケーションを切り替えるサーバを決定する仕組み（スプリットブレイン解決）がある。このスプリットブレイン解決方式では、ネットワーク機器に障害が発生したサーバが共有ディスクへアクセスを行うと、障害サーバが健全に作動していると誤認識するおそれがあり、障害サーバから確実に共有ディスクにアクセスできないようにする必要がある。

【0006】

上記のように、クラスタシステムにおいてフェイルオーバを行う場合には、正しく排他制御を行なって障害サーバから共有ディスクへのアクセスを確実に停止する必要がある。

【先行技術文献】

【特許文献】

10

【0007】

【特許文献1】特開2012-173752号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

ここで、クラスタシステムにおいて共有ディスクへのアクセスを停止する手法として、ディスクのアンマウントやFC（Fibre Channel）スイッチのポート閉塞、OS（Operating System）パニック、などが用いられている。ところが、上述した各手法には以下の様な課題がある。

【0009】

20

ディスクのアンマウントの場合は、処理に時間がかかり、書き込み中のプロセスが存在する場合にはアンマウント処理が失敗する、という問題がある。FCスイッチのポート閉塞の場合は、サーバ外のモジュールに接続するため接続時間がかかり、障害の種類によってはFCスイッチに接続できない、という問題がある。OSパニックの場合は、HBA（Host Bus Adapter）カードのキャッシュに残っているI/Oが書き込まれる可能性がある、という問題が生じる。

【0010】

また、HW（Hardware）レイヤで共有ディスクへの書き込みを防ぐI/Oフェンシングを実施する場合は、カードへの電源供給を断つ方式などがある。HWレイヤでI/Oフェンシングを実施した場合は、確実なI/Oフェンシングが可能であるが、I/Oフェンシングの完了を待機系サーバに通知する手段がない。このため、例えば、待機系サーバは、ハートビートの通信タイムアウトを待ってフェイルオーバを開始しなければならず、迅速にフェイルオーバを行うことができない、という問題がある。

30

【0011】

以上のように、フェイルオーバ型のクラスタシステムでは、共有データの保護を行いつつ、迅速にフェイルオーバを行うことができない、という問題があった。

【0012】

このため、本発明の目的は、上述した課題である、共有データの保護を行いつつ、迅速にフェイルオーバを行うことができない、ということを解決することができるクラスタシステムを提供することにある。

40

【課題を解決するための手段】

【0013】

本発明の一形態であるクラスタシステムは、相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバ機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備え、前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えており、

前記現用系サーバの前記制御装置は、前記待機系サーバと通信を行う通信モジュールと

50

、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールと、を備え、

前記待機系サーバは、前記現用系サーバの前記制御装置から前記I/Oフェンシングが完了した旨の通知を受けて、前記待機系サーバ上における前記アプリケーションにて前記現用系サーバ上における前記アプリケーションによる処理を引き継ぐフェイルオーバー処理機能を有する、
という構成をとる。

【0014】

また、本発明の他の形態である現用系サーバは、

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムにおける前記現用系サーバであって、

前記現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えており、

前記制御装置は、前記待機系サーバと通信を行う通信モジュールと、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールと、を備えた、
という構成をとる。

【0015】

また、本発明の他の形態であるプログラムは、

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムにおける前記現用系サーバが、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えと共に、
当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記制御装置に、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールによる処理、
を実行させるためのプログラムである。

【0016】

また、本発明の他の形態であるフェイルオーバー方法は、

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムによるフェイルオーバー方法であって、

前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えと共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記現用系サーバの前記制御装置は、前記現用系サーバの障害発生時に、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行い、当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバ

10

20

30

40

50

に通知し、

前記待機系サーバは、前記現用系サーバの前記制御装置から前記I/Oフェンシングが完了した旨の通知を受けて、前記待機系サーバ上における前記アプリケーションにて前記現用系サーバ上における前記アプリケーションによる処理を引き継ぐ、という構成をとる。

【0017】

また、本発明の他の形態であるフェイルオーバー方法は、

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えた

10

クラスタシステムによるフェイルオーバー方法であって、
前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記現用系サーバの前記制御装置は、前記現用系サーバの障害発生時に、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行い、当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、

という構成をとる。

20

【発明の効果】

【0018】

本発明は、以上のように構成されることにより、共有データの保護を行いつつ、迅速にフェイルオーバーを行うことができる。

【図面の簡単な説明】

【0019】

【図1】本発明の実施形態1におけるクラスタシステムの構成を示すブロック図である。

【図2】図1に開示したクラスタシステムにおけるクラスタ起動時の動作を示すフローチャートである。

【図3】図1に開示したクラスタシステムにおけるクラスタ起動時の動作を説明する図である。

30

【図4】図1に開示したクラスタシステムにおけるソフトウェア障害検出時の動作を示すフローチャートである。

【図5】図1に開示したクラスタシステムにおけるソフトウェア障害検出時の動作を説明する図である。

【図6】図1に開示したクラスタシステムにおけるハードウェア障害検出時の動作を示すフローチャートである。

【図7】図1に開示したクラスタシステムにおけるハードウェア障害検出時の動作を説明する図である。

【図8】図1に開示したクラスタシステムにおける待機系サーバによる現用系サーバの障害検出時の動作を示すフローチャートである。

40

【図9】図1に開示したクラスタシステムにおける待機系サーバによる現用系サーバの障害検出時の動作を説明する図である。

【図10】本発明の付記1におけるクラスタシステムの構成を示すブロック図である。

【発明を実施するための形態】

【0020】

<実施形態1>

本発明の第1の実施形態を、図1乃至図9を参照して説明する。図1は、クラスタシステムの構成を示す図であり、図2乃至図9は、クラスタシステムの動作を示す図である。

【0021】

50

〔構成〕

本実施形態におけるクラスタシステム1は、共有ディスクを備えたフェイルオーバ型のクラスタシステムである。このクラスタシステム1は、図1に示すように、現用系サーバ10と待機系サーバ20とを備えており、これら両サーバ10, 20にてアクセス可能であり共有される共有ディスク50を備えて構成されている。

【0022】

現用系サーバ10と待機系サーバ20とは、それぞれほぼ同一の構成を取っており、同一の業務アプリケーションを実行することが可能である。そして、クラスタシステム1は、現用系サーバ10に障害が生じた場合には、当該現用系サーバ10の業務アプリケーションによる処理を、待機系サーバ20の業務アプリケーションが引き継ぐフェイルオーバ機能を有している。以下、現用系サーバ10と待機系サーバ20の構成はほぼ同一であるため、主に現用系サーバ10の構成について詳述する。

10

【0023】

現用系サーバ10は、演算装置と記憶装置とを備えた情報処理装置にて構成されている。そして、現用系サーバ10は、演算装置にプログラムが組み込まれることで構築されたOS(Operating System)20を備えており、当該OS20は、クラスタソフトウェア30とIPMIドライバ21とを備えている。また、OS20は、図示しない業務アプリケーションを備えている。

【0024】

また、現用系サーバ10は、当該サーバ10に装備されたハードウェア(HW)を管理するBMC(Base Management Controller)40(制御装置)と、共有ディスク50へのアクセスを行うPCIカード11(ディスク入出力装置)を備えている。なお、上記BMC40は、上記OS20の影響を受けずに作動するよう構成されている。以下、各構成についてさらに詳述する。

20

【0025】

上記OS20が備えるクラスタソフトウェア30は、サーバ管理モジュール31、HB(Heart Beat)モジュール32、障害検出モジュール33、サーバ制御モジュール34、通知受信モジュール35を備えている。

【0026】

上記サーバ管理モジュール31は、クラスタシステムを構成する各サーバの状態(起動状態/停止状態)や、業務アプリケーションの状態(起動状態/停止状態、起動しているサーバ)を管理する。また、サーバ管理モジュール31は、HBモジュール32や障害検出モジュール33から障害情報や他サーバの停止情報を受信した場合に、業務アプリケーションの状態を確認し、アプリケーションが停止している、または、正常に動作できないと判断した場合に、業務アプリケーションのフェイルオーバを指示する。

30

【0027】

上記HBモジュール32は、定期的なネットワークパケットの送受信やディスクへの書込などの手段によりクラスタシステムを構成する各サーバが起動しているか停止しているかを判断する。他のサーバが起動状態から停止状態、または、停止状態から起動状態に移行した場合は、サーバ管理モジュール31にサーバの状態変化を通知する。

40

【0028】

上記障害検出モジュール33は、自サーバの内蔵・共有ディスクへのアクセス可否や業務アプリケーションの死活、ネットワーク通信が正常に行えるかなど、自サーバのコンピュータ資源が正常に動作しているかどうかの監視を行う。コンピュータ資源に異常を検出した場合は、サーバ管理モジュール31に異常を通知する。

【0029】

上記サーバ制御モジュール34は、サーバ管理モジュール31からの指示に従い自サーバ10のシャットダウン、自サーバ10のBMC40や相手サーバ20のBMC90への指示を行い、自サーバ10や相手サーバ20の起動停止を制御する。

【0030】

50

通知受信モジュール35は、相手サーバ20のBMC90からのI/Oフェンシング完了通知を受信し、サーバ管理モジュール31に通知を行う。

【0031】

そして、上述したOS20及びクラスタソフトウェア30の構成を、待機系サーバ60のOS70及びクラスタソフトウェア80も備えている。つまり、待機系サーバ60のOS70は、IPMIドライバ71を備えており、クラスタソフトウェア80は、サーバ管理モジュール81、HB(Heart Beat)モジュール82、障害検出モジュール83、サーバ制御モジュール84、通知受信モジュール85を備えている。

【0032】

次に、現用系サーバ10のBMC40の構成について詳述する。BMC40は、当該BMC40が備える演算装置にプログラムが組み込まれることで構築された、BMC制御モジュール41、IPMI受信モジュール42、フェンシング制御モジュール43、通知送信モジュール44、HW障害検出モジュール45を備えている。

【0033】

BMC制御モジュール41(初期化モジュール)は、起動時にクラスタソフトウェア30から通知されるI/Oフェンシング対象となるPCIカード11の-slot番号や、フェンシング完了通知先となる待機系サーバ20のIPアドレスの管理を行う。なお、PCIカード11の-slot番号は、PCIカード11にアクセスするためのアクセス情報であり、待機系サーバ20のIPアドレスは、情報を通知する際に当該待機系サーバ20の宛先となるアドレス情報である。

【0034】

また、BMC制御モジュール41は、クラスタソフトウェア30からI/Oフェンシング命令を受信した場合や、HW障害検出モジュール45からHW障害の通知を受信した場合に、I/Oフェンシングの実行をフェンシング制御モジュール43に指示する。

【0035】

上記IPMI受信モジュール42は、自サーバ10のクラスタソフトウェア30や他サーバ20のクラスタソフトウェア80からのIPMI制御コマンド(I/Oフェンシング命令や設定情報)を受信し、BMC制御モジュール41に通知する。

【0036】

上記フェンシング制御モジュール43(初期化モジュール)は、BMC制御モジュール41からのI/Oフェンシング命令を受信し、指定されたPCIカード11が接続されているPCI-slotのリセットを行い、I/Oフェンシングを行う。また、フェンシング制御モジュール43は、I/Oフェンシング完了後に、通知送信モジュール44にI/Oフェンシングの完了を通知する。なお、上記I/Oフェンシングとは、クラスタシステムを構成するサーバが停止する場合に、サーバをFile System(共有ディスク)から切り離してI/Oが発生しない状態にすることである。

【0037】

上記通知送信モジュール44(通信モジュール)は、I/Oフェンシングの完了を受信すると、登録されたIPアドレスに対して、I/Oフェンシングの完了をネットワーク経由で通知する。

【0038】

上記HW障害検出モジュール45は、サーバ10を構成するHWモジュールを監視し、障害検出時にBMC制御モジュール41に障害を通知する。

【0039】

そして、上述したBMC40の構成を、待機系サーバ60のBMC90も備えている。つまり、待機系サーバ60のBMC90は、当該BMC90が備える演算装置にプログラムが組み込まれることで構築された、BMC制御モジュール91、IPMI受信モジュール92、フェンシング制御モジュール93、通知送信モジュール94、HW障害検出モジュール95を備えている。

【0040】

10

20

30

40

50

[動作]

次に、上述したクラスタシステムの動作を、図 2 乃至図 9 を参照して説明する。以下では、クラスタシステムの動作を、次の 4 種類のケースに分けて説明する。

- (1) クラスタ起動
- (2) クラスタソフトウェアが自サーバの障害を検出したケース
- (3) BMC が自サーバの HW 障害を検出したケース
- (4) クラスタソフトウェアが、他サーバの停止を検出したケース

【 0 0 4 1 】

- (1) クラスタ起動

まず、クラスタシステムが起動したときの動作を、図 2 のフローチャート、及び、図 3 の点線矢印を参照して説明する。

【 0 0 4 2 】

まず、現用系サーバ 1 0 の起動時に、サーバ制御モジュール 3 4 は、予め登録された設定ファイルから、共有ディスク 5 0 と接続している PCI カード 1 1 のスロット番号と、待機系サーバ 2 0 の IP アドレスと、をロードし (ステップ S 1)、IPMI ドライバ 2 1 経由で BMC 4 0 に通知する (ステップ S 2)。通知を受信した BMC 4 0 の IPMI 受信モジュール 4 2 は、BMC 制御モジュール 4 1 に PCI カード 1 1 のスロット番号と待機系サーバの IP アドレスを通知する。そして、BMC 制御モジュール 4 1 は、通知を受けた PCI カード 1 1 のスロット番号と待機系サーバ 2 0 の IP アドレスとを、モジュール内の揮発性メモリに保存する (ステップ S 3)。なお、待機系サーバ 2 0 でも同様の処理が行われる。

【 0 0 4 3 】

上記の初期設定が行われない場合は、HB タイムアウトによる他サーバの停止検出と、OS パニックやディスクのアンマウントにより共有ディスクへの I/O を停止する一般的なクラスタソフトウェアとして動作する。

【 0 0 4 4 】

- (2) クラスタソフトウェアが自サーバの障害を検出したケース

次に、クラスタソフトウェアが自サーバの障害を検出したときの動作を、図 4 のフローチャート、及び、図 5 の点線矢印を参照して説明する。本障害ケースでは、現用系サーバ 1 0 上のクラスタソフトウェア 3 0 を構成する障害検出モジュール 3 3 が障害を検出する。本ケースにおける障害は、業務アプリケーションの消滅、NIC ドライバやディスクドライバなどのドライバレベルで検出される障害など、ソフトウェアレベルで検出可能な障害である。

【 0 0 4 5 】

まず、障害検出モジュール 3 3 が障害を検出すると、サーバ管理モジュール 3 1 に障害の発生を通知する (ステップ S 1 1)。サーバ管理モジュール 3 1 は、BMC 4 0 による I/O フェンシング機能が有効であることを確認し、速やかに業務アプリケーションのフェイルオーバー処理を開始する。

【 0 0 4 6 】

サーバ管理モジュール 3 1 は、共有ディスク 5 0 への I/O の停止を行うために、サーバ制御モジュール 3 4 に共有ディスク 5 0 への I/O フェンシング処理を指示する (ステップ S 1 2)。サーバ制御モジュール 3 4 は、サーバ管理モジュール 3 1 からの I/O フェンシング指示を受け取ると直ちに、I/O フェンシング命令と対象の PCI カード 1 1 のスロット番号を、IPMI ドライバ 2 1 経由で BMC 4 0 に通知する。

【 0 0 4 7 】

続いて、IPMI 受信モジュール 4 2 は、IPMI ドライバ 2 1 経由で I/O フェンシング命令と対象 PCI カードのスロット番号を受信する。IPMI 受信モジュール 4 2 は、BMC 制御モジュール 4 1 に I/O フェンシング命令と対象となる PCI カード 1 1 のスロット番号を通知する。BMC 制御モジュール 4 1 は、フェンシング制御モジュール 4 3 に I/O フェンシング命令と対象となる PCI カード 1 1 のスロット番号を通知する。なお、BMC 制御モジュール 4 1 は、起動時に登録された PCI カード 1 1 のスロット番号を、フェンシング制御モジュール 4 3 に通

10

20

30

40

50

知してもよい。

【 0 0 4 8 】

そして、フェンシング制御モジュール 4 3 は、I/Oフェンシングの命令を受けると、指定されたPCIカード 1 1 のスロット番号に対して当該PCIカード 1 1 をリセット（初期化）する（ステップ S 1 3）。これにより、ソフトウェアレベルでI/Oフェンシングを制御するよりも、確実に共有ディスク 5 0 に対するI/Oを停止することができる。

【 0 0 4 9 】

上述したように、PCIカード 1 1 のリセットが完了すると、フェンシング制御モジュール 4 3 は、通知送信モジュール 4 4 にI/Oフェンシングの完了を通知するように指示を行う。通知送信モジュール 4 4 は、あらかじめ登録されている通知先（待機系サーバ 2 0 のIPアドレス）に対して、I/Oフェンシングの完了通知を送信する（ステップ S 1 4）。通知の方法は、ネットワーク経由のSNMP Trapなどを用いる。なお、上述したBMC 4 0 がI/Oフェンシング命令を受信してから完了通知を待機系サーバ 2 0 に送信する一連の処理を、BMCによるI/Oフェンシング処理と呼ぶ。

10

【 0 0 5 0 】

その後、待機系サーバ 2 0 上のクラスタウェア 8 0 を構成する通知受信モジュール 8 5 は、ネットワーク経由で、現用系サーバ 1 0 のBMC 4 0 から当該現用系サーバ 1 0 のI/Oフェンシングの完了通知を受信し、I/Oフェンシングの完了したサーバの情報をサーバ管理モジュール 8 1 に通知する。サーバ管理モジュール 8 1 は、I/Oフェンシングの完了通知と対象サーバである現用系サーバ 1 0 の情報を受信すると、当該サーバ 1 0 で動作していた業務アプリケーションのフェイルオーバを開始する。これにより、業務アプリケーションのフェイルオーバが完了し、待機系サーバ 2 0 にて業務が再開する（ステップ S 1 5）。

20

【 0 0 5 1 】

ここで、クラスタシステムを構成するサーバが 3 台以上の場合は、業務アプリケーションのフェイルオーバ先を、予め登録された優先順位や、リソースの空き状況などのルールに従って 1 台のサーバに決定する。なお、待機系サーバ 2 0 の通知受信モジュール 8 5 がI/Oフェンシングの完了通知を受信してからフェイルオーバを行う一連の処理を、クラスタ再構成処理と呼ぶ。

【 0 0 5 2 】

（ 3 ） BMCが自サーバのHW障害を検出したケース

次に、現用系サーバ 1 0 のBMC 4 0 が、自サーバのHW障害を検出したときの動作を、図 6 のフローチャート、及び、図 7 の点線矢印を参照して説明する。

30

【 0 0 5 3 】

本ケースでは、現用系サーバ 1 0 のBMC 4 0 を構成するHW障害検出モジュール 4 5 が障害を検出する（ステップ S 2 1）。本ケースにおける障害は、ハードウェア（HW）のセンサ異常や温度上昇、電圧低下などのハードウェアレベルで検出可能な障害である。

【 0 0 5 4 】

まず、HW障害検出モジュール 4 5 は、障害を検出した場合、BMC制御モジュール 4 1 に障害の発生を通知する。BMC制御モジュール 4 1 は、障害内容を確認し、現用系サーバ 1 0 の停止が必要な障害と判断した場合は、フェンシング制御モジュール 4 3 にI/Oフェンシング命令と起動処理で保存したPCIカード 1 1 のスロット番号を、フェンシング制御モジュール 4 3 に通知する（ステップ S 2 2）。

40

【 0 0 5 5 】

フェンシング制御モジュール 4 3 は、指定された番号のPCIカード 1 1 のスロットをリセットする（ステップ S 2 2）。これにより、共有ディスク 5 0 に対するI/Oを停止するI/Oフェンシングを行う。

【 0 0 5 6 】

PCIカード 1 1 のスロットのリセットが完了すると、フェンシング制御モジュール 4 3 は通知送信モジュール 4 4 に、I/Oフェンシングの完了を通知するように指示を行う。通

50

知送信モジュール 44 は、あらかじめ登録されている通知先（待機系サーバ 20 の IP アドレス）に対して、I/O フェンシングの完了通知を送信する（ステップ S 23）。

【0057】

待機系サーバ 20 のクラスタウェア 80 は、I/O フェンシングの完了通知を受信し、クラスタ再構成処理を行う（ステップ S 24）。

【0058】

（4）クラスタソフトウェアが、他サーバの停止を検出したケース

次に、待機系サーバ 60 のクラスタソフトウェア 80 が、現用系サーバの障害を検出したときの動作を、図 8 のフローチャート、及び、図 9 の点線矢印を参照して説明する。

【0059】

本ケースでは、待機系サーバ 20 のクラスタソフトウェア 80 を構成する HB モジュール 82 が、現用系サーバ 10 の障害を検出する。本ケースにおける障害は、現用系サーバ 10 の OS 20 のストールやパニック、クラスタソフトウェア 30 の停止である。

【0060】

まず、HB モジュール 82 は、現用系サーバ 10 の HB モジュール 32 と定期的にハートビート通信を行なっている。このとき、HB モジュール 82 が一定期間ハートビート通信ができない状態になると、当該 HB モジュール 82 は現用系サーバ 10 が停止したと判断し（ステップ S 31）、サーバ管理モジュール 81 にハートビート通信異常を通知する。サーバ管理モジュール 81 は、スプリットブレイン解決を行い、現用系サーバ 10 が停止していると判断した場合に、サーバ制御モジュール 84 に現用系サーバ 10 の I/O フェンシング実行を指示する。

【0061】

本処理は、現用系サーバ 10 の OS 20 がストールによりハートビート通信できていない場合、共有ディスク 50 に対する I/O が完全に停止していないことになるため、確実な I/O フェンシングを実施するための処理である。サーバ制御モジュール 84 は、サーバ管理モジュール 81 からの I/O フェンシング命令を受信すると直ちに、ネットワーク経由で I/O フェンシング命令と対象の PCI カード 11 のスロット番号を、現用系サーバ 10 の BMC 40 に通知する（ステップ S 32）。

【0062】

現用系サーバ 10 の BMC 40 の IPMI 受信モジュール 42 は、待機系サーバ 60 のサーバ制御モジュール 84 からの I/O フェンシング命令を受信すると、BMC 40 による I/O フェンシング処理を行う。つまり、上述同様に、共有ディスク 50 に接続された PCI カード 11 のスロットをリセットし（ステップ S 33）、あらかじめ登録されている通知先（待機系サーバ 20 の IP アドレス）に対して、I/O フェンシングの完了通知を送信する（ステップ S 34）。

【0063】

その後、待機系サーバ 60 のクラスタウェア 80 は、I/O フェンシングの完了通知を受信し、クラスタ再構成処理を行う（ステップ S 35）。

【0064】

本発明のクラスタシステムは以上のように作動するため、以下の効果を発揮しうる。

まず、共有ディスクアクセスの排他制御の確実性が向上する。これは、上述したように、BMC および HW レベルで PCI カードをリセットするためであり、ソフトウェアレベルで I/O フェンシングを制御するよりも確実に I/O を停止可能である。たとえば、ソフトウェアレベルでは、カードにキューイングされている I/O を障害発生時に停止することは困難である。また、OS がストールまたは停止していると判断される障害サーバに対して、外部のサーバが I/O フェンシングの命令を行うことで、確実に I/O フェンシングを実行できる。

【0065】

また、高速なフェイルオーバー実現できる。これは、I/O フェンシング完了後に、現用系サーバの BMC から待機系サーバのクラスタウェアにネットワーク通知を行っていることによる。通常のクラスタウェアは、ハートビートの通信異常（タイムアウト）を検知して現

10

20

30

40

50

用系サーバが停止したと判断するため、必ずタイムアウト時間が必要となる。本発明では、タイムアウトではなく、BMCからの通知を契機にフェイルオーバを実行するため、フェイルオーバ開始までの時間短縮が可能となる。

【0066】

ここで、上記では、2台のクラスタ構成であったが、I/Oフェンシング通知対象IPアドレスを複数登録することで、3台以上のクラスタ構成でも同様の方式が実現可能である。また、両サーバアクティブのクラスタ構成でもアクティブ-スタンバイ型と同用の構成で、提案の方式が実現可能である。

【0067】

<付記>

上記実施形態の一部又は全部は、以下の付記のようにも記載されうる。以下、本発明におけるクラスタシステム（図10参照）、現用系サーバ、プログラム、フェイルオーバ方法の構成の概略を説明する。但し、本発明は、以下の構成に限定されない。

【0068】

(付記1)

相互に同一のアプリケーション111, 211を実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバ機能を有する現用系サーバ100及び待機系サーバ200と、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスク300と、を備え、

前記現用系サーバ100は、当該現用系サーバに組み込まれたオペレーティングシステム110の影響を受けずに作動する制御装置120と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置130と、を備えており、

前記現用系サーバ100の前記制御装置120は、前記待機系サーバと通信を行う通信モジュール121と、前記現用系サーバの障害発生時に前記ディスク入出力装置130を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュール121を介して前記待機系サーバ200に通知する初期化モジュール122と、を備え、

前記待機系サーバ200は、前記現用系サーバの前記制御装置から前記I/Oフェンシングが完了した旨の通知を受けて、前記待機系サーバ上における前記アプリケーションにて前記現用系サーバ上における前記アプリケーションによる処理を引き継ぐフェイルオーバ処理機能を有する、クラスタシステム。

【0069】

(付記2)

付記1に記載のクラスタシステムであって、

前記現用系サーバは、当該現用系サーバにて作動しているソフトウェアの障害発生を検出して前記制御装置の前記通知モジュールに通知する障害検出モジュールを備え、

前記制御装置の前記初期化モジュールは、前記障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、クラスタシステム。

【0070】

(付記3)

付記1又は2に記載のクラスタシステムであって、

前記現用系サーバの前記制御装置は、当該現用系サーバに装備されたハードウェアの障害発生を検出して前記制御装置の前記初期化モジュールに通知するハードウェア障害検出モジュールを備え、

前記制御装置の前記初期化モジュールは、前記ハードウェア障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化す

10

20

30

40

50

ることにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、クラスタシステム。

【0071】

(付記4)

付記1乃至3のいずれかに記載のクラスタシステムであって、

前記待機系サーバは、前記現用系サーバにおける障害発生を検出して、前記現用系サーバの前記初期化モジュールに通知する待機系側障害検出モジュールを備え、

前記現用系サーバの前記初期化モジュールは、前記待機系サーバの前記待機系側障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、クラスタシステム。

【0072】

(付記5)

付記1乃至4のいずれかに記載のクラスタシステムであって、

前記初期化モジュールは、前記ディスク入出力装置にアクセスするためのアクセス情報と、前記待機系サーバに情報を通知するためのアドレス情報と、を記憶保持しており、前記アクセス情報に基づいて前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記アドレス情報に基づいて前記通信モジュールを介して前記待機系サーバに通知する、クラスタシステム。

【0073】

(付記6)

付記5に記載のクラスタシステムであって、

前記現用系サーバは、当該現用系サーバの起動時に、前記アクセス情報及び前記アドレス情報が予め設定された設定ファイルを読み出し、当該設定ファイルに設定されている前記アクセス情報及び前記アドレス情報を前記制御装置の前記初期化モジュールに通知し、

前記初期化モジュールは、通知された前記アクセス情報及び前記アドレス情報を記憶保持する、クラスタシステム。

【0074】

(付記7)

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムにおける前記現用系サーバであって、

前記現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えており、

前記制御装置は、前記待機系サーバと通信を行う通信モジュールと、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールと、を備えた、現用系サーバ。

【0075】

(付記7-2)

付記7に記載の現用系サーバであって、

前記現用系サーバは、当該現用系サーバにて作動しているソフトウェアの障害発生を検

10

20

30

40

50

出して前記初期化モジュールに通知する障害検出モジュールを備え、

前記初期化モジュールは、前記障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
現用系サーバ。

【0076】

(付記7-3)

付記7又は7-2に記載の現用系サーバであって、

前記制御装置は、当該現用系サーバに装備されたハードウェアの障害発生を検出して前記初期化モジュールに通知するハードウェア障害検出モジュールを備え、

前記初期化モジュールは、前記ハードウェア障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
現用系サーバ。

【0077】

(付記7-4)

付記7乃至7-3のいずれかに記載の現用系サーバであって、

前記待機系サーバが、前記現用系サーバにおける障害発生を検出して、前記現用系サーバの前記初期化モジュールに通知する待機系側障害検出モジュールを備えており、

前記初期化モジュールは、前記待機系サーバの前記待機系側障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
現用系サーバ。

【0078】

(付記8)

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムにおける前記現用系サーバが、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記制御装置に、前記現用系サーバの障害発生時に前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する初期化モジュールによる処理、

を実行させるためのプログラム。

【0079】

(付記8-2)

付記8に記載のプログラムであって、

前記現用系サーバは、当該現用系サーバにて作動しているソフトウェアの障害発生を検出して前記初期化モジュールに通知する障害検出モジュールを備えており、

前記初期化モジュールは、前記障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対するI/Oフェンシングを行うと共に当該I/Oフェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
プログラム。

【 0 0 8 0 】

(付記 8 - 3)

付記 8 又は 8 - 2 に記載のプログラムであって、

前記制御装置は、当該現用系サーバに装備されたハードウェアの障害発生を検出して前記初期化モジュールに通知するハードウェア障害検出モジュールを備えており、

前記初期化モジュールは、前記ハードウェア障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
プログラム。

10

【 0 0 8 1 】

(付記 8 - 4)

付記 8 乃至 8 - 3 のいずれかに記載のプログラムであって、

前記待機系サーバは、前記現用系サーバにおける障害発生を検出して、前記現用系サーバの前記初期化モジュールに通知する待機系側障害検出モジュールを備えており、

前記現用系サーバの前記初期化モジュールは、前記待機系サーバの前記待機系側障害検出モジュールから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
プログラム。

20

【 0 0 8 2 】

(付記 9)

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムによるフェイルオーバー方法であって、

前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

30

前記現用系サーバの前記制御装置は、前記現用系サーバの障害発生時に、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行い、当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知し、

前記待機系サーバは、前記現用系サーバの前記制御装置から前記 I/O フェンシングが完了した旨の通知を受けて、前記待機系サーバ上における前記アプリケーションにて前記現用系サーバ上における前記アプリケーションによる処理を引き継ぐ、
フェイルオーバー方法。

【 0 0 8 3 】

40

(付記 9 - 2)

付記 9 に記載のフェイルオーバー方法であって、

前記現用系サーバは、当該現用系サーバにて作動しているソフトウェアの障害発生を検出して前記制御装置に通知し、

前記現用系サーバの前記制御装置は、前記現用系サーバにて作動しているソフトウェアの障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
フェイルオーバー方法。

【 0 0 8 4 】

50

(付記 9 - 3)

付記 9 又は 9 - 2 に記載のフェイルオーバー方法であって、

前記現用系サーバの前記制御装置は、当該現用系サーバに装備されたハードウェアの障害発生を検出すると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
フェイルオーバー方法。

【0085】

(付記 9 - 4)

付記 9 乃至 9 - 3 のいずれかに記載のフェイルオーバー方法であって、

前記待機系サーバは、前記現用系サーバにおける障害発生を検出して、前記現用系サーバの前記制御装置に通知し、

前記現用系サーバの前記制御装置は、前記待機系サーバから前記現用系サーバにおける障害発生の通知を受けると、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行うと共に当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
フェイルオーバー方法。

【0086】

(付記 10)

相互に同一のアプリケーションを実行し当該アプリケーションによる他方の処理を引き継ぐことが可能なフェイルオーバー機能を有する現用系サーバ及び待機系サーバと、当該現用系サーバ及び待機系サーバにて共有可能なデータを記憶する共有ディスクと、を備えたクラスタシステムによるフェイルオーバー方法であって、

前記現用系サーバは、当該現用系サーバに組み込まれたオペレーティングシステムの影響を受けずに作動する制御装置と、前記共有ディスクに接続し当該共有ディスクに対するデータの入出力を行うディスク入出力装置と、を備えると共に、当該制御装置は、前記待機系サーバと通信を行う通信モジュールを備えており、

前記現用系サーバの前記制御装置は、前記現用系サーバの障害発生時に、前記ディスク入出力装置を初期化することにより前記共有ディスクに対する I/O フェンシングを行い、当該 I/O フェンシングが完了した旨を前記通信モジュールを介して前記待機系サーバに通知する、
フェイルオーバー方法。

【0087】

なお、上述したプログラムは、記憶装置に記憶されていたり、コンピュータが読み取り可能な記録媒体に記録されている。例えば、記録媒体は、フレキシブルディスク、光ディスク、光磁気ディスク、及び、半導体メモリ等の可搬性を有する媒体である。

【0088】

以上、上記実施形態等を参照して本願発明を説明したが、本願発明は、上述した実施形態に限定されるものではない。本願発明の構成や詳細には、本願発明の範囲内で当業者が理解しうる様々な変更をすることができる。

【符号の説明】

【0089】

- 10 現用系サーバ
- 11 PCIカード
- 20 OS
- 30 クラスタソフトウェア
- 31 サーバ管理モジュール
- 32 HBモジュール
- 33 障害検出モジュール
- 34 サーバ制御モジュール

10

20

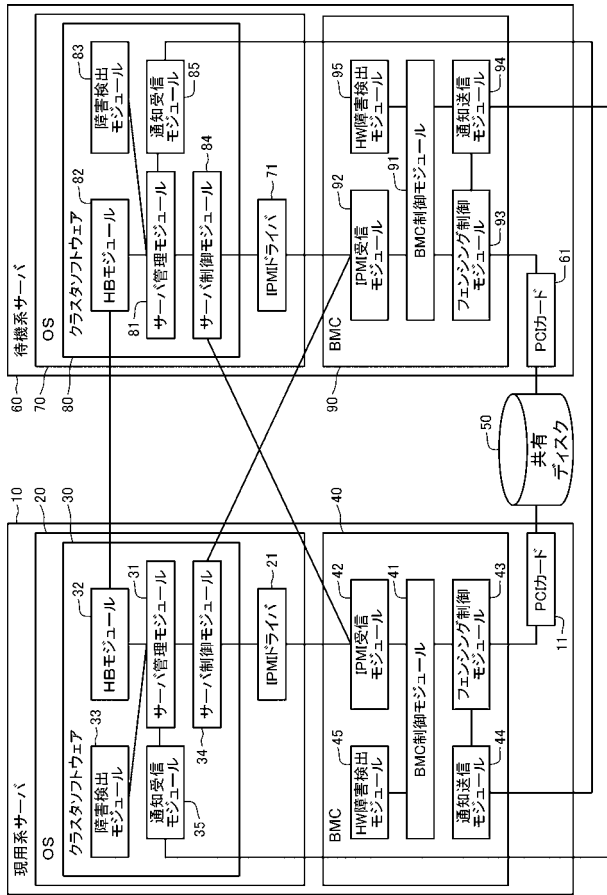
30

40

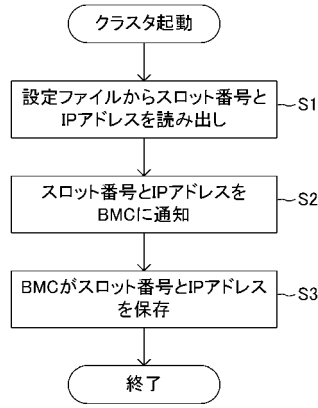
50

3 5	通知受信モジュール	
4 0	BMC	
4 1	BMC制御モジュール	
4 2	IPMI受信モジュール	
4 3	フェンシング制御モジュール	
4 4	通知送信モジュール	
4 5	HW障害検出モジュール	
5 0	共有ディスク	
6 0	待機系サーバ	
6 1	PCIカード	10
7 0	OS	
8 0	クラスタソフトウェア	
8 1	サーバ管理モジュール	
8 2	HBモジュール	
8 3	障害検出モジュール	
8 4	サーバ制御モジュール	
8 5	通知受信モジュール	
9 0	BMC	
9 1	BMC制御モジュール	
9 2	IPMI受信モジュール	20
9 3	フェンシング制御モジュール	
9 4	通知送信モジュール	
9 5	HW障害検出モジュール	
1 0 0	現用系サーバ	
1 1 0	OS	
1 1 1	アプリケーション	
1 2 0	制御装置	
1 2 1	通信モジュール	
1 2 2	初期化モジュール	
1 3 0	入出力装置	30
2 0 0	待機系サーバ	
2 1 0	OS	
2 1 1	アプリケーション	
3 0 0	共有ディスク	

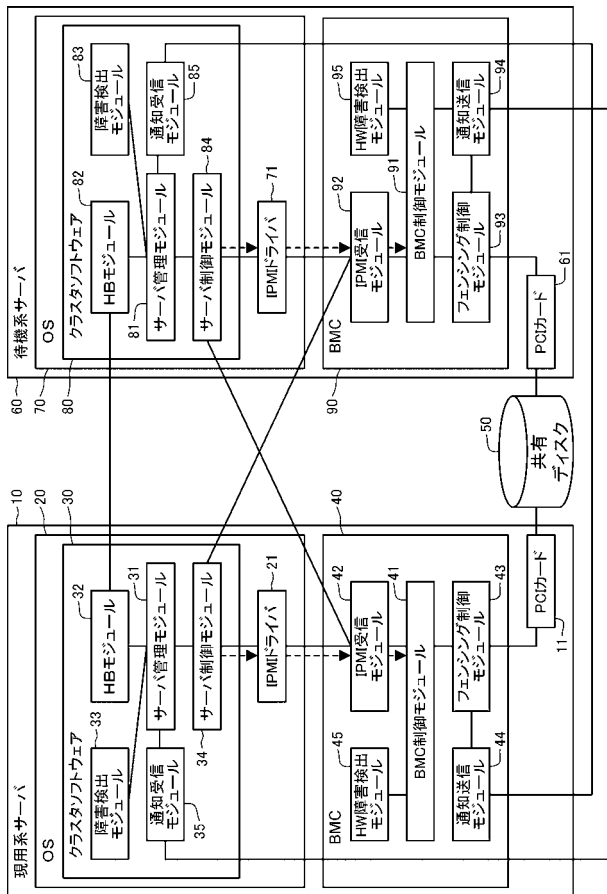
【 図 1 】



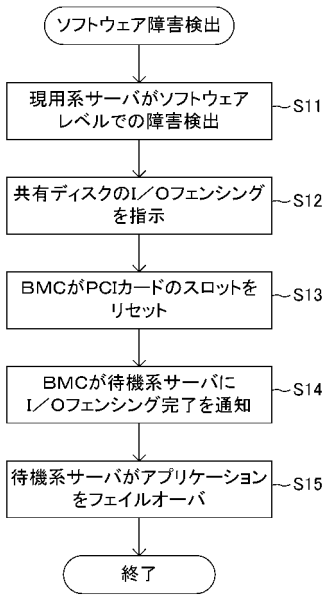
【 図 2 】



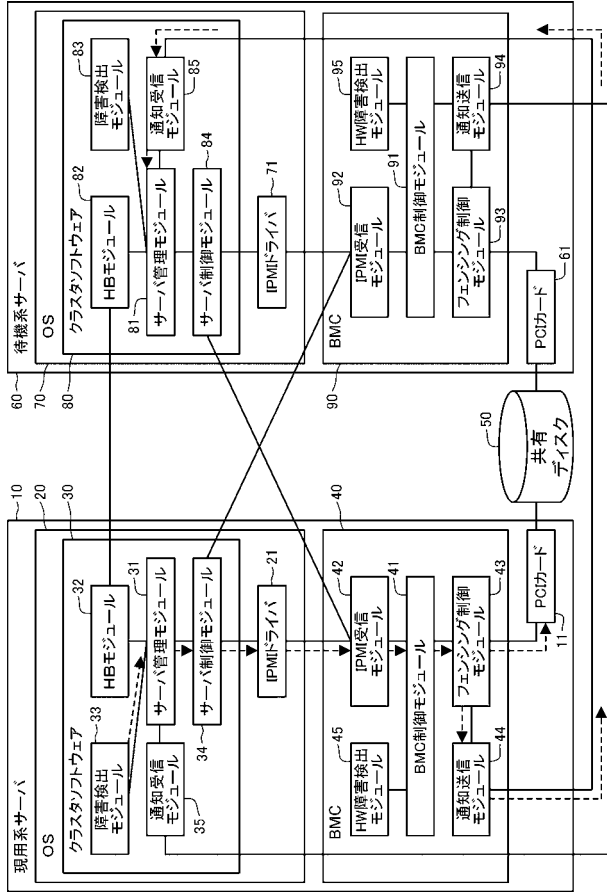
【 図 3 】



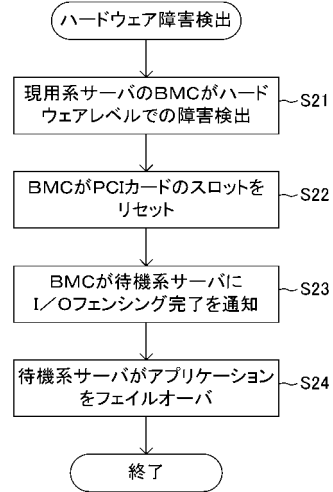
【 図 4 】



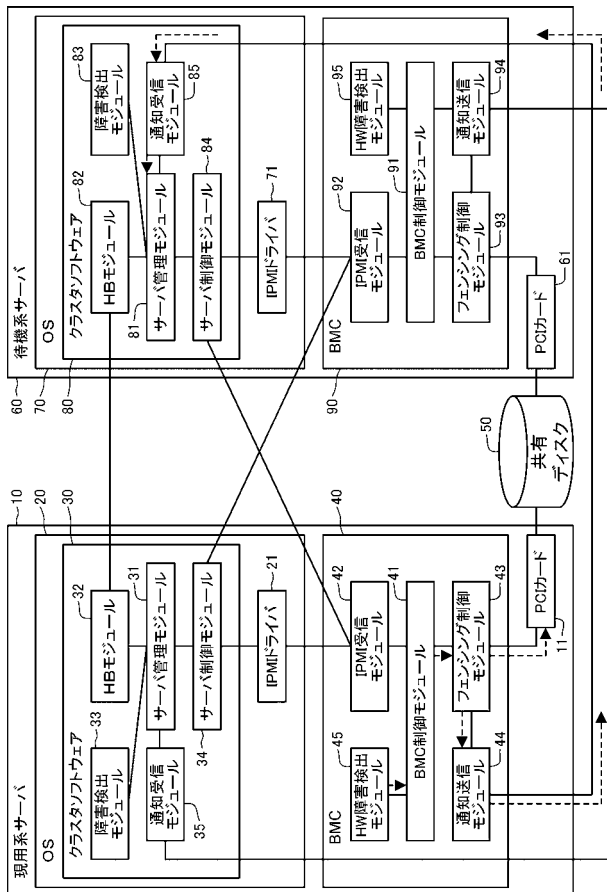
【図5】



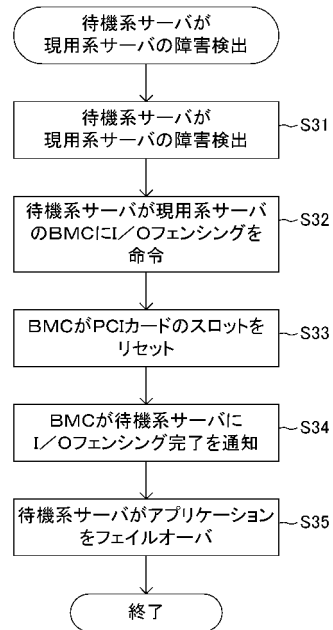
【図6】



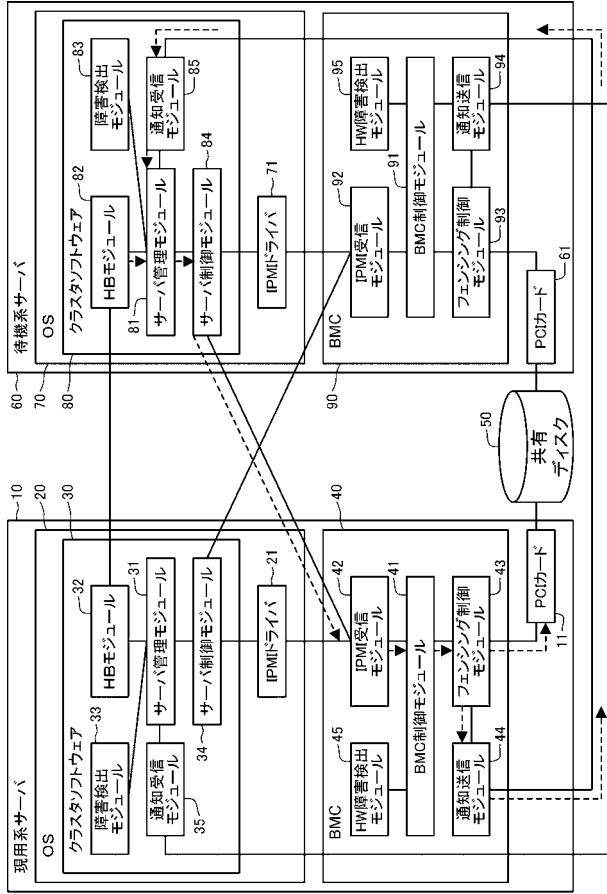
【図7】



【図8】



【図9】



【図10】

