

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6070064号  
(P6070064)

(45) 発行日 平成29年2月1日(2017.2.1)

(24) 登録日 平成29年1月13日(2017.1.13)

(51) Int. Cl. F 1  
**G 0 6 F** 11/20 (2006.01) G 0 6 F 11/20 6 9 7  
**G 0 6 F** 9/50 (2006.01) G 0 6 F 9/46 4 6 5

請求項の数 5 (全 18 頁)

(21) 出願番号 特願2012-237731 (P2012-237731)  
 (22) 出願日 平成24年10月29日(2012.10.29)  
 (65) 公開番号 特開2014-89506 (P2014-89506A)  
 (43) 公開日 平成26年5月15日(2014.5.15)  
 審査請求日 平成27年9月9日(2015.9.9)

(73) 特許権者 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 100079164  
 弁理士 高橋 勇  
 (72) 発明者 森内 哲  
 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 三坂 敏夫

最終頁に続く

(54) 【発明の名称】 ノード装置、クラスタシステム、フェイルオーバー方法およびプログラム

(57) 【特許請求の範囲】

【請求項1】

他のノード装置と相互に接続されてクラスタシステムを構成すると共に、前記他のノード装置との間で共有される外部記憶装置である共有ストレージに接続されてなるノード装置であって、

メモリと、

接続された各クライアント装置のIPアドレスと当該クライアント装置のための前記メモリとの間の対応関係を予め記憶しているIPアドレステーブル記憶手段と、

前記クライアント装置からの要求に基づいて予め装備されたアプリケーションソフトを実行して処理を行い、これによって得られる処理データを処理依頼元の前記クライアント装置のIPアドレスに対応する前記メモリに記憶させるアプリケーションソフト実行部と

10

、記憶された前記処理データを前記他のノード装置のメモリに記憶させるミラーリング処理部と、

前記他のノード装置のメモリに記憶させられた後の前記処理データを前記共有ストレージに書き込むストレージ記憶部と、

前記他のノード装置のメモリに前記処理データが残った状態で当該他のノード装置に異常が発生した場合に、前記IPアドレステーブル記憶手段に記憶された対応関係を変更して当該他のノード装置による処理を引き継ぐフェイルオーバー処理部と、

を有することを特徴とするノード装置。

20

## 【請求項 2】

前記ミラーリング処理部が、前記処理データを前記他のノード装置のメモリに記憶させる処理の完了後、処理要求元の前記クライアント装置に書き込み終了通知を返信する機能を有することを特徴とする、請求項 1 に記載のノード装置。

## 【請求項 3】

第 1 および第 2 のノード装置と、前記第 1 および第 2 のノード装置の間で共有される外部記憶装置である共有ストレージとが相互に接続されて構築されたクラスタシステムであって、

前記第 1 および第 2 のノード装置が、請求項 1 又は 2 に記載のノード装置であることを特徴とするクラスタシステム。

10

## 【請求項 4】

他のノード装置と相互に接続されてクラスタシステムを構成すると共に、前記他のノード装置との間で共有される外部記憶装置である共有ストレージに接続されてなるノード装置にあって、

接続された各クライアント装置の IP アドレスと当該クライアント装置のためのメモリとの間の対応関係が予め備えられた IP アドレステーブル記憶手段に記憶されたものであると共に、

アプリケーションソフトをアプリケーションソフト実行部が実行し、

前記アプリケーションソフトによって得られた処理データを前記クライアント装置の IP アドレスに対応する前記メモリ上に前記アプリケーションソフト実行部が一時的に保存し、

20

記憶された前記処理データをミラーリング処理部が前記他のノード装置の対応する記憶領域に記憶させ、

前記他のノード装置のメモリに記憶させられた後の前記処理データをストレージ記憶部が前記共有ストレージに書き込み、

前記他のノード装置のメモリに前記処理データが残った状態で当該他のノード装置に異常が発生した場合に、前記 IP アドレステーブル記憶手段に記憶された対応関係をフェイルオーバー処理部が変更して当該他のノード装置による処理を引き継ぐ、

ことを特徴とするフェイルオーバー方法。

## 【請求項 5】

30

他のノード装置と相互に接続されてクラスタシステムを構成すると共に、前記他のノード装置との間で共有される外部記憶装置である共有ストレージに接続されてなるノード装置にあって、

接続された各クライアント装置の IP アドレスと当該クライアント装置のためのメモリとの間の対応関係が予め備えられた IP アドレステーブル記憶手段に記憶されたものであると共に、

前記ノード装置が備えるプロセッサに、

アプリケーションソフトを実行する手順、

前記アプリケーションソフトによって得られた処理データを前記クライアント装置の IP アドレスに対応する前記メモリ上に一時的に保存する手順、

40

記憶された前記処理データを前記他のノード装置の対応するメモリに記憶させる手順、

前記他のノード装置のメモリに記憶させられた後の前記処理データを前記共有ストレージに書き込む手順、

および前記他のノード装置のメモリに前記処理データが残った状態で当該他のノード装置に異常が発生した場合に、前記 IP アドレステーブル記憶手段に記憶された対応関係をフェイルオーバー処理部が変更して当該他のノード装置による処理を引き継ぐ手順、

を実行させることを特徴とするフェイルオーバープログラム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

50

本発明はノード装置、クラスタシステム、フェイルオーバー方法およびプログラムに関し、特に異常発生時に短時間にフェイルオーバー処理を可能とするノード装置等に関する。

【背景技術】

【0002】

企業などで利用されるコンピュータシステムにおいては、短時間の停止であっても、その間に発生した業務の停止によって巨額の損失が発生することがある。そのため、そのようなコンピュータシステムにおいては「高可用性(High Availability)」、より具体的には「稼働率99.99%以上(年間停止時間52分以下)」が求められている。そのようなコンピュータシステムでは、サーバクラスタと呼ばれる構成が多く利用されている。

10

【0003】

サーバクラスタは、同一のNAS(Network Attached Storage、ネットワーク接続ストレージ)を参照する複数台のサーバコンピュータが1台の仮想サーバとして動作するように構成したものである。サーバクラスタは、このNASの利用形態により、「アクティブ-アクティブ構成」と「アクティブ-パッシブ構成」の2種類に分類される。

【0004】

アクティブ-アクティブ構成のサーバクラスタは、当該サーバクラスタを構成する各サーバがいずれも稼働系サーバとして動作して処理を行う。このため、オフロード(off-load、負荷分散)も兼ねることができ、システム全体でのパフォーマンスを向上させることも可能である。

20

【0005】

図10は、特許文献1に記載されている、既存のクラスタシステム901の構成について示す説明図である。クラスタシステム901は、第一ノード910および第二ノード920という各々のコンピュータ装置(サーバ)と、その他多数のクライアント装置940とが、ネットワーク941を介して相互に接続されて構成されている。

【0006】

第一ノード910および第二ノード920は、共有の外部記憶装置である共有ストレージ930と接続されている。そして、第一ノード910および第二ノード920は、処理にかかる負荷を相互に分散しつつ、記憶されているデータを互いにミラーリングして、どちらか一方に故障が発生した場合にその故障した方の装置で行われていた処理を残る一方が引き継いで続行する(これをフェイルオーバーという)ことができる構成となっている。

30

【0007】

第一ノード910は、コンピュータプログラムを実行する主体であるプロセッサ911と、処理中のデータを一時的に記憶する不揮発性の主記憶装置であるNVRAM(Non-Volatile RAM)912と、処理されたデータを共有ストレージ930に固定的に記憶する外部ストレージ接続手段913と、ネットワーク941を介して第二ノード920や各クライアント装置940との間で通信を行う通信手段914とが備えられている。

【0008】

プロセッサ911は、クライアント装置940からの依頼に基づく処理を行うアプリケーションソフトを実行してそのデータをNVRAM912に書き込むアプリケーション実行部951と、NVRAM912との間のデータ交換を仲介するデバイスドライバを動作させるNVRAMドライバ実行部952と、外部ストレージ接続手段913を経由して共有ストレージ930との間のデータ交換を仲介するデバイスドライバを動作させるストレージ装置ドライバ実行部953と、通信手段914との間のデータ交換を仲介するデバイスドライバを動作させる外部通信ドライバ954として機能する。

40

【0009】

NVRAM912は、第一ノード910で行われる処理についてデータを一時的に書き込む第一ノード用領域912aと、第二ノード920で行われる処理で後述のNVRAM922に書き込まれたデータをミラーリングする第二ノード用領域912bとに分かれて

50

いる。第一ノード910では、クライアント装置940からの依頼に応じて第一ノード用領域912aに対してデータ処理を行い、その内容を第二ノード920側の第一ノード用領域922aにもコピーして反映させる。

【0010】

プロセッサ911はさらに、各々のプログラムの動作により、アプリケーション実行部951がクライアント装置940からの依頼に応じて第一ノード用領域912aに書き込んだデータを第二ノード920の第一ノード用領域922aにコピーして反映させるミラーリング処理部955と、そのデータを共有ストレージ930に記憶するストレージ記憶部956と、第二ノード920で異常が発生した場合にその処理を引き継ぐフェイルオーバー処理部957としても機能する。

10

【0011】

第二ノード920は、ハードウェア的にもソフトウェア的にも、次に説明する点を除いては第一ノード910と同一の構成を有するので、呼称は全て同一とし、参照番号は各々+10ずつという。即ち、ハードウェアとしてはプロセッサ921、NVRAM922...などのようにいい、ソフトウェアとしてはアプリケーション実行部961、NVRAMドライバ実行部962...などのようにいう。

【0012】

第二ノード920の、第一ノード910と比べての唯一の相違点について説明する。NVRAM922は、第一ノード910で行われる処理でNVRAM912に書き込まれたデータをミラーリングする第一ノード用領域922aと、第二ノード920で行われる処理についてデータを一時的に書き込む第二ノード用領域922bとに分かれている。第二ノード920では、クライアント装置940からの依頼に応じて第二ノード用領域922bに対してデータ処理を行い、その内容を第一ノード910側の第二ノード用領域912bにもコピーして反映させる。

20

【0013】

第一ノード910のアプリケーション実行部951がデータをNVRAM912の第一ノード用領域912aに書き込む際には、既にファイルシステムに書き込み可能な状態で記憶させる。第二ノード920のアプリケーション実行部961がデータをNVRAM922の第二ノード用領域922bに書き込む際も同様である。この状態のデータに対して、第一ノード910および第二ノード920は互いにミラーリング処理部955および965によってNVRAM912および922のデータを相互にミラーリングし、そしてストレージ記憶部956もしくは966によって共有ストレージ930に記憶させる。

30

【0014】

ここで、たとえば第一ノード910で、NVRAM912上に共有ストレージ930に未反映のデータが残った状態で異常が発生した場合には、第二ノード920側でフェイルオーバー処理部967が、第一ノード用領域922a上のデータに対してリカバリ処理、即ち共有ストレージ930に記憶させる処理を行った後で、各クライアント装置940からのアクセス受付を再開させる。第二ノード920で同様の異常が発生した場合には、第一ノード910がこれと同様の動作を行う。

【0015】

このように構成することで、第一ノード910と第二ノード920とで各クライアント装置940からのアクセスにかかる処理の負荷を分散しつつ、一方で異常が発生した時には残る一方ですぐに処理を再開することが可能となる。この構造は、自ノード用の領域を使用する限りにおいては、シングルシステムと同様の構造で動作することができる。このため、クラスタ構成として特別な改造が必要なく、通常運用での構造変更を必要としないでクラスタシステムの構築が可能となる。

40

【0016】

他にこれに関連する技術文献として、たとえば次の各文献がある。その中でも特許文献2には、NVRAMを複数ブロックに分割して各ブロックごとにリード/ライトプロテクトすることが可能であるというディスク制御装置について記載されている。特許文献3に

50

は、遠隔地に設置されたリモートサイトにデータをミラーリングするシステムについて記載されている。非特許文献1および2には、NASクラスタモデルの基本的な構成について記載されている。

【先行技術文献】

【特許文献】

【0017】

【特許文献1】特開2007-328778号公報

【特許文献2】特開2008-217811号公報

【特許文献3】特表2006-527875号公報

【非特許文献】

【0018】

【非特許文献1】デジタルアドバンテージ「第20回 ファイル共有プロトコルSMB/CIFS(その1)(基礎から学ぶWindowsネットワーク Windowsネットワーク管理者への道 より)」、平成16年10月29日、アイティメディア株式会社、[平成24年5月21日検索]、インターネット<URL: [http://www.atmarkit.co.jp/fwin2k/network/baswinlan020/baswinlan020\\_01.html](http://www.atmarkit.co.jp/fwin2k/network/baswinlan020/baswinlan020_01.html)>

【非特許文献2】高橋郷「Windowsクラスタリング入門 第1回 MSCS導入の準備～サーバ・クラスタの基礎知識～」、平成20年12月3日、アイティメディア株式会社、[平成23年12月21日検索]、インターネット<URL: [http://www.atmarkit.co.jp/fwin2k/operation/mscluster01/mscluster01\\_01.html](http://www.atmarkit.co.jp/fwin2k/operation/mscluster01/mscluster01_01.html)>

【発明の概要】

【発明が解決しようとする課題】

【0019】

図10で説明した特許文献1記載の既存のクラスタシステム901は、前述したように、通常動作時には負荷を分散しつつ、異常発生時にもすぐに処理を再開することが可能である。

【0020】

しかしながらこの構成では、NVRAM912または922上のデータを共有ストレージ930に反映を完了させないと、データの整合性を保つことができない。この処理の途中で異常が発生した場合、共有ストレージ930に途中まで書き込んでいたデータの整合性は失われてしまう。

【0021】

前述の第一ノード910で異常が発生した場合の例でいえば、第二ノード920側で第一ノード用領域922a上のデータに対するリカバリ処理、即ち当該データを共有ストレージ930に記憶させる処理が、第二ノード920が第一ノード910の動作を引き継ぐためには必要である。第二ノード920で異常が発生した場合も同様である。

【0022】

従って、そのリカバリ処理を行うための時間が必要であり、クライアント装置940に異常の発生を報告する必要が生じることとなる。この問題を解決しうる技術は、残る特許文献2～3および非特許文献1～2にも記載されていない。

【0023】

本発明の目的は、クライアント装置側に異常の発生を意識させることなく、ごく短時間でフェイルオーバー処理を行うことを可能とするクラスタシステム、ノード、フェイルオーバー方法およびプログラムを提供することにある。

【課題を解決するための手段】

【0024】

上記目的を達成するため、本発明に係るノード装置は、他のノード装置と相互に接続されてクラスタシステムを構成すると共に、他のノード装置との間で共有される外部記憶装置である共有ストレージに接続されてなるノード装置であって、メモリと、接続された各クライアント装置のIPアドレスと当該クライアント装置のためのメモリとの間の対応関

10

20

30

40

50

係を予め記憶しているIPアドレステーブル記憶手段と、クライアント装置からの要求に基づいて予め装備されたアプリケーションソフトを実行して処理を行い、これによって得られる処理データを処理依頼元のクライアント装置のIPアドレスにメモリに記憶させるアプリケーションソフト実行部と、記憶された処理データを他のノード装置のメモリに記憶させるミラーリング処理部と、他のノード装置のメモリに記憶させられた後の処理データを共有ストレージに書き込むストレージ記憶部と、他のノード装置のメモリに処理データが残った状態で当該他のノード装置に異常が発生した場合に、IPアドレステーブル記憶手段に記憶された対応関係を変更して当該他のノード装置による処理を引き継ぐフェイルオーバー処理部と、を有することを特徴とする。

【0025】

上記目的を達成するため、本発明に係るクラスタシステムは、第1および第2のノード装置と、第1および第2のノード装置の間で共有される外部記憶装置である共有ストレージとが相互に接続されて構築されたクラスタシステムであって、第1および第2のノード装置が、請求項1又は2に記載のノード装置であることを特徴とする。

【0026】

上記目的を達成するため、本発明に係るフェイルオーバー方法は、他のノード装置と相互に接続されてクラスタシステムを構成すると共に、他のノード装置との間で共有される外部記憶装置である共有ストレージに接続されてなるノード装置にあって、接続された各クライアント装置のIPアドレスと当該クライアント装置のためのメモリとの間の対応関係が予め備えられたIPアドレステーブル記憶手段に記憶されたものであると共に、アプリケーションソフトをアプリケーションソフト実行部が実行し、アプリケーションソフトによって得られた処理データをクライアント装置のIPアドレスに対応するメモリ上にアプリケーションソフト実行部が一時的に保存し、記憶された処理データをミラーリング処理部が他のノード装置の対応する記憶領域に記憶させ、他のノード装置のメモリに記憶させられた後の処理データをストレージ記憶部が共有ストレージに書き込み、他のノード装置のメモリに処理データが残った状態で当該他のノード装置に異常が発生した場合に、IPアドレステーブル記憶手段に記憶された対応関係をフェイルオーバー処理部が変更して当該他のノード装置による処理を引き継ぐことを特徴とする。

【0027】

上記目的を達成するため、本発明に係るフェイルオーバープログラムは、他のノード装置と相互に接続されてクラスタシステムを構成すると共に、他のノード装置との間で共有される外部記憶装置である共有ストレージに接続されてなるノード装置にあって、接続された各クライアント装置のIPアドレスと当該クライアント装置のためのメモリとの間の対応関係が予め備えられたIPアドレステーブル記憶手段に記憶されたものであると共に、ノード装置が備えるプロセッサに、アプリケーションソフトを実行する手順、アプリケーションソフトによって得られた処理データをクライアント装置のIPアドレスに対応するメモリ上に一時的に保存する手順、記憶された処理データを他のノード装置の対応するメモリに記憶させる手順、他のノード装置の対応する記憶領域に記憶させられた後の処理データを共有ストレージに書き込む手順、および他のノード装置のメモリに処理データが残った状態で当該他のノード装置に異常が発生した場合に、IPアドレステーブル記憶手段に記憶された対応関係をフェイルオーバー処理部が変更して当該他のノード装置による処理を引き継ぐ手順を実行させることを特徴とする。

【発明の効果】

【0028】

本発明は、上記したように、各クライアント装置のIPアドレスに対応して不揮発性メモリの記憶領域を区切って、ミラーリング処理を行ってから共有ストレージに書き込むように構成したので、一方のノード装置で異常が発生しても不揮発性メモリに記憶されたデータの整合性は保たれる。これによってクライアント装置側に異常の発生を意識させることなく、ごく短時間でフェイルオーバー処理を行うことが可能であるという、優れた特徴を持つクラスタシステム、ノード、フェイルオーバー方法およびプログラムを提供するこ

10

20

30

40

50

とができる。

【図面の簡単な説明】

【0029】

【図1】本発明の実施形態に係るクラスタシステムの構成について示す説明図である。

【図2】図1に示した第一ノードに備えられるIPアドレステーブル記憶手段の記憶内容について示す説明図である。

【図3】図1に示したクラスタシステムで、第一ノードおよび第二ノードの正常動作時の処理の分担について示す説明図である。

【図4】図1に示したクラスタシステムで、第一ノードおよび第二ノードの異常発生時の処理の分担について示す説明図である。

【図5】図1に示したクラスタシステムの正常時の動作について示す説明図である。

【図6】図1に示したクラスタシステムの正常時の動作について示すフローチャートである。

【図7】図6の続きである。

【図8】図1に示したクラスタシステムの異常発生時の動作について示す説明図である。

【図9】図1に示したクラスタシステムの異常発生時の動作について示すフローチャートである。

【図10】特許文献1に記載されている、既存のクラスタシステムの構成について示す説明図である。

【発明を実施するための形態】

【0030】

(実施形態)

以下、本発明の実施形態の構成について添付図1に基づいて説明する。

最初に、本実施形態の基本的な内容について説明し、その後でより具体的な内容について説明する。

本実施形態に係るノード装置(第一ノード10および第二ノード20)は、同一の構成を有する他のノード装置と相互に接続されてクラスタシステムを構成すると共に、他のノード装置との間で共有される同一の外部記憶装置である共有ストレージに接続されてなるノード装置である。このノード装置(第一ノード10)は、予め複数の記憶領域に区切られた不揮発性メモリ(NVRAM12)と、接続された各クライアント装置40のIPアドレスと当該クライアント装置が使用すべき不揮発性メモリの記憶領域との間の対応関係を予め記憶しているIPアドレステーブル記憶手段15と、クライアント装置からの要求に基づいて予め装備されたアプリケーションソフトを実行して処理を行い、これによって得られる処理データを処理依頼元のクライアント装置のIPアドレスに対応する記憶領域上に記憶させるアプリケーションソフト実行部101と、記憶された処理データを他のノード装置の対応する記憶領域に記憶させるミラーリング処理部105と、他のノード装置の対応する記憶領域に記憶させられた後の処理データを共有ストレージに書き込むストレージ記憶部106とを有する。

【0031】

ここで、ミラーリング処理部105は、処理データを他のノード装置の対応する記憶領域に記憶させる処理の完了後、処理要求元のクライアント装置40に書き込み終了通知を返信する機能を有する。さらに、他のノード装置の不揮発性メモリに処理データが残った状態で当該他のノード装置に異常が発生した場合に、IPアドレステーブル記憶手段に記憶された対応関係を変更して当該他のノード装置による処理を引き継ぐフェイルオーバー処理部107も備える。

【0032】

以上の構成を備えることにより、ノード装置(第一ノード10)は、クライアント装置側に異常の発生を意識させることなく、ごく短時間でフェイルオーバー処理を行うことが可能となる。

以下、これをより詳細に説明する。

## 【 0 0 3 3 】

図 1 は、本発明の実施形態に係るクラスタシステム 1 の構成について示す説明図である。クラスタシステム 1 は、第一ノード 1 0 および第二ノード 2 0 という各々のコンピュータ装置（サーバ）と、その他多数のクライアント装置 4 0 とが、ネットワーク 4 1 を介して相互に接続されて構成されている。

## 【 0 0 3 4 】

第一ノード 1 0 および第二ノード 2 0 は、単にクライアント装置 4 0 から受信したデータを保存するファイルサーバでもよいし、また受信したデータに対して何らかの処理を行ってから保存するデータベースシステム、ウェブサーバ、業務システムなどでもよい。

## 【 0 0 3 5 】

またここで、第一ノード 1 0 および第二ノード 2 0 とクライアント装置 4 0 との間で利用される通信方式は、クライアント装置 4 0 の OS（Operating System: 基本ソフト）がたとえばウィンドウズ（登録商標）であれば C I F S（Common Internet File System）を利用することができるし、U N I X（登録商標）であれば N I S（Network File System）を利用することができる。これら以外にも、利用可能な通信方式であれば任意のものを使用することができる。

## 【 0 0 3 6 】

ここでサーバクラスタを構成している第一ノード 1 0 および第二ノード 2 0 は、共有の外部記憶装置である共有ストレージ 3 0 と接続されている。そして、第一ノード 1 0 および第二ノード 2 0 は、処理にかかる負荷を相互に分散しつつ、記憶されているデータを互いにミラーリングして、どちらか一方に故障が発生した場合にその故障した方の装置で行われていた処理を残る一方が引き継いで続行する（これをフェイルオーバーという）ことができる構成となっている。

## 【 0 0 3 7 】

第一ノード 1 0 は、コンピュータプログラムを実行する主体であるプロセッサ 1 1 と、処理中のデータを一時的に記憶する不揮発性の主記憶装置である N V R A M 1 2 と、処理されたデータを共有ストレージ 3 0 に固定的に記憶する外部ストレージ接続手段 1 3 と、ネットワーク 4 1 を介して第二ノード 2 0 や各クライアント装置 4 0 との間で通信を行う通信手段 1 4 と、後述の I P アドレステーブル記憶手段 1 5 とが備えられている。

## 【 0 0 3 8 】

N V R A M 1 2 は、I C A（InterConnectAccess）カードとして第一ノード 1 0 内に実装されている。この方式によって接続されることにより、N V R A M 1 2 はプロセッサ 1 1 上で動作するオペレーティングシステムを経由することなく、R D M A（Remote Direct Memory Access）プロトコルによってデータを第二ノード 2 0 に転送することが可能となるであるというメリットがある。他にも、N V R A M 1 2 を実装する上では、任意の接続方式を利用することができる。

## 【 0 0 3 9 】

プロセッサ 1 1 は、クライアント装置 4 0 からの依頼に基づく処理を行うアプリケーションソフトを実行してそのデータを N V R A M 1 2 に書き込むアプリケーション実行部 1 0 1 と、N V R A M 1 2 との間のデータ交換を仲介するデバイスドライバを動作させる N V R A M ドライバ実行部 1 0 2 と、外部ストレージ接続手段 1 3 を経由して共有ストレージ 3 0 との間のデータ交換を仲介するデバイスドライバを動作させるストレージ装置ドライバ実行部 1 0 3 と、通信手段 1 4 との間のデータ交換を仲介するデバイスドライバを動作させる外部通信ドライバ実行部 1 0 4 として機能する。

## 【 0 0 4 0 】

N V R A M 1 2 には、第一ノード 1 0 および第二ノード 2 0 で各々行われる処理について、後述するようにクライアント装置 4 0 の I P アドレスのグループに応じてグループ 0 用領域 1 2 a およびグループ 1 用領域 1 2 b とに分かれている。第一ノード 1 0 では、「グループ 0」に属するクライアント装置 4 0 からの依頼に応じてグループ 0 用領域 1 2 a に対してデータ処理を行い、その内容を第二ノード 2 0 側の同領域にもコピーして反映す

10

20

30

40

50

る。

【0041】

プロセッサ11はさらに、各々のプログラムの動作により、アプリケーション実行部101が「グループ0」に属するクライアント装置40からの依頼に応じてグループ0用領域12aに書き込んだデータを第二ノード20のグループ0用領域22aにコピーして反映させるミラーリング処理部105と、そのデータを共有ストレージ30に記憶するストレージ記憶部106と、第二ノード20で異常が発生した場合にその処理を引き継ぐようIPアドレステーブル記憶手段15の記憶内容を変更するフェイルオーバー処理部107としても機能する。

【0042】

第二ノード20は、第一ノード10と、ハードウェア的にもソフトウェア的にも同一の構成を有している。従って、第二ノード20の各機能部については、第一ノード10と呼称を全て同一とし、参照番号は最上位の「1」を「2」に代えた以外は全て第一ノード10と同一とする。即ち、ハードウェア的にはプロセッサ21、NVRAM22...などのようにいい、ソフトウェア的にはアプリケーション実行部201、NVRAMドライバ実行部202...などのようにいう。

【0043】

図2は、図1に示した第一ノード10に備えられるIPアドレステーブル記憶手段15の記憶内容について示す説明図である。クライアント装置40は、各々のIPアドレスのグループに応じて「グループ0」および「グループ1」の2グループに分かれる。IPアドレステーブル記憶手段15には、各クライアント装置40のIPアドレスと、それに該当するクライアント装置40が「グループ0」および「グループ1」のうちのいずれに属するかについて記憶されている。第二ノード20にも、これと同一の内容を記憶するIPアドレステーブル記憶手段25が存在する。

【0044】

ここでいう「グループ0」および「グループ1」の分け方については、たとえばIPアドレスの範囲によって分けてもよいし、1個ごとのIPアドレスについて「グループ0」もしくは「グループ1」を指定してもよい。また、IPアドレスではなくホスト名で「グループ0」もしくは「グループ1」に分類してもよい。

【0045】

図3は、図1に示したクラスタシステム1で、第一ノード10および第二ノード20の正常動作時の処理の分担について示す説明図である。第一ノード10のNVRAM12は前述のようにグループ0用領域12aおよびグループ1用領域12bとに分かれ、第二ノード20のNVRAM22も同様にグループ0用領域22aおよびグループ1用領域22bとに分かれる。

【0046】

IPアドレステーブル記憶手段15(25)に記憶された各クライアント装置40の「グループ0」および「グループ1」の区分について、第一ノード10のアプリケーション実行部101は、「グループ0」に属するクライアント装置40からの依頼に応じて処理を行い、その処理結果をNVRAM12のグループ0用領域12aに書き込む。第二ノード20のアプリケーション実行部201は、「グループ1」に属するクライアント装置40からの依頼に応じて処理を行い、その処理結果をNVRAM22のグループ1用領域22bに書き込む。

【0047】

そして、「グループ0」に属するクライアント装置40からの依頼に応じた処理の場合、第一ノード10のアプリケーション実行部101がデータ処理を行った後で第一ノード10および第二ノード20のミラーリング処理部105および205は、NVRAM12および22の相互の内容をコピーしあい、そして第一ノード10のストレージ記憶部106がそのデータをストレージ装置ドライバ実行部103および外部ストレージ接続手段13を介して共有ストレージ30にそのデータを記憶する。

10

20

30

40

50

## 【 0 0 4 8 】

「グループ 1」に属するクライアント装置 4 0 からの依頼に応じた処理の場合は、第二ノード 2 0 のアプリケーション実行部 2 0 1 が先に処理を行う点以外は、「グループ 0」の場合と同一である。

## 【 0 0 4 9 】

図 4 は、図 1 に示したクラスタシステム 1 で、第一ノード 1 0 および第二ノード 2 0 の異常発生時の処理の分担について示す説明図である。第一ノード 1 0 の側で、共有ストレージ 3 0 への書き込みの済んでいないデータが N V R A M 1 2 に残っている状態で第一ノード 1 0 が異常を起こして停止した場合、ミラーリング処理部 1 0 5 によって、N V R A M 2 2 のグループ 0 用領域 2 2 a にその未反映データがコピーされている。

10

## 【 0 0 5 0 】

そこで、第二ノード 2 0 が「グループ 0」に属するクライアント装置 4 0 からの依頼によるその処理を引き継ぎ、クライアント装置 4 0 からの依頼に応じた処理によるデータを、N V R A M 2 2 のグループ 0 用領域 2 2 a にコピーされた分のデータに続いて書き込んで、これを共有ストレージ 3 0 に記憶する。

## 【 0 0 5 1 】

図 5 は、図 1 に示したクラスタシステム 1 の正常時の動作について示す説明図である。図 6 ~ 7 ( 図面の錯綜回避のため 2 枚に分ける ) は、図 1 に示したクラスタシステム 1 の正常時の動作について示すフローチャートである。より詳しくは、図 5 では、「グループ 0」に属するクライアント装置 4 0 からの処理依頼に応じての動作を示している。

20

## 【 0 0 5 2 】

クライアント装置 4 0 からデータ処理依頼がなされた場合、第一ノード 1 0 および第二ノード 2 0 のアプリケーション実行部 1 0 1 および 2 0 1 はそれぞれ、各々の I P アドレステーブル記憶手段 1 5 および 2 5 を参照して、「グループ 0」および「グループ 1」のどちらに属するクライアント装置 4 0 からの処理依頼かを判断する ( ステップ S 3 0 1 および 3 5 1 ) 。

## 【 0 0 5 3 】

「グループ 0」に属するクライアント装置 4 0 からの処理依頼だった場合、第一ノード 1 0 のアプリケーション実行部 1 0 1 がその依頼に対応する処理を行い、その処理結果を N V R A M 1 2 のグループ 0 用領域 1 2 a に書き込む ( ステップ S 3 0 2 ) 。そしてミラーリング処理部 1 0 5 がそのデータを第二ノード 2 0 のミラーリング処理部 2 0 5 に渡して記憶させる ( ステップ S 3 0 3 ) 。これを受けた第二ノード 2 0 のミラーリング処理部 2 0 5 は、そのデータを N V R A M 2 2 のグループ 0 用領域 2 2 a に書き込む ( ステップ S 3 0 7 ) 。

30

## 【 0 0 5 4 】

そして、ミラーリング処理部 1 0 5 が依頼元のクライアント装置 4 0 に書き込み終了通知を返し ( ステップ S 3 0 4 ) 、ストレージ記憶部 1 0 6 がクライアント装置 4 0 に対して書き込み終了を通知するタイミングでそのデータをストレージ装置ドライバ実行部 1 0 3 および外部ストレージ接続手段 1 3 を介して共有ストレージ 3 0 にそのデータを記憶する ( ステップ S 3 0 5 ) 。共有ストレージ 3 0 への書き込みが終了したら、ミラーリング処理部 1 0 5 はそのデータを N V R A M 1 2 のグループ 0 用領域 1 2 a から削除する ( ステップ S 3 0 6 ) 。

40

## 【 0 0 5 5 】

「グループ 1」に属するクライアント装置 4 0 からの処理依頼だった場合、第二ノード 2 0 のアプリケーション実行部 2 0 1 がその依頼に対応する処理を行い、その処理結果を N V R A M 2 2 のグループ 1 用領域 2 2 b に書き込む ( ステップ S 3 5 2 ) 。そしてミラーリング処理部 2 0 5 がそのデータを第一ノード 1 0 のミラーリング処理部 1 0 5 に渡して記憶させる ( ステップ S 3 5 3 ) 。これを受けた第一ノード 1 0 のミラーリング処理部 1 0 5 は、そのデータを N V R A M 1 2 のグループ 1 用領域 1 2 b に書き込む ( ステップ S 3 5 7 ) 。

50

## 【 0 0 5 6 】

そして、ミラーリング処理部 2 0 5 が依頼元のクライアント装置 4 0 に書き込み終了通知を返し（ステップ S 3 5 4）、ストレージ記憶部 2 0 6 がそのデータをストレージ装置ドライバ実行部 2 0 3 および外部ストレージ接続手段 2 3 を介して共有ストレージ 3 0 にそのデータを記憶する（ステップ S 3 5 5）。共有ストレージ 3 0 への書き込みが終了したら、ミラーリング処理部 2 0 5 はそのデータを N V R A M 2 2 のグループ 1 用領域 2 2 a から削除する（ステップ S 3 5 6）。

## 【 0 0 5 7 】

図 8 は、図 1 に示したクラスタシステム 1 の異常発生時の動作について示す説明図である。図 9 は、図 1 に示したクラスタシステム 1 の異常発生時の動作について示すフローチャートである。より詳しくは、「グループ 0」に属するクライアント装置 4 0 からの処理依頼に応じて第一ノード 1 0 が図 5 ・ステップ S 3 0 5 に示した共有ストレージ 3 0 への書き込み処理を行っている間に、この第一ノード 1 0 に故障が発生して、N V R A M 1 2 のグループ 0 用領域 1 2 a に共有ストレージ 3 0 へ未反映のデータが残ってしまった場合について、図 8 ~ 9 では示している。

10

## 【 0 0 5 8 】

この段階であれば、ミラーリング処理部 1 0 5 によって、N V R A M 2 2 のグループ 0 用領域 2 2 a にその未反映データがコピーされている。そこで、異常発生を検出したフェイルオーバー処理部 2 0 7 は（ステップ S 4 0 1）、第二ノード 2 0 の I P アドレステーブル記憶手段 2 5 で、「グループ 0」に属するクライアント装置 4 0 の処理を第二ノード 2 0 が引き継ぐように設定し直す（ステップ S 4 0 2）。

20

## 【 0 0 5 9 】

ここで、第一ノード 1 0 および第二ノード 2 0 では、互いが正常に動作していることを確認するために、常時周期的にハートビート通信を行っている。本実施形態で「N V R A M に未反映データが残ったままの状態、一方のノードで異常が発生した」ことを検出して上記ステップ S 4 0 1 からの動作開始は、このハートビート通信に対する返答が一定時間ないことを検出した場合をその契機とすることができる。

## 【 0 0 6 0 】

これによって、N V R A M 2 2 のグループ 0 用領域 2 2 a に残っている未反映データの続きの処理をアプリケーション実行部 2 0 1 が行ってそのデータを N V R A M 2 2 のグループ 0 用領域 2 2 a に書き込み（ステップ S 4 0 3）、ミラーリング処理部 2 0 5 が依頼元のクライアント装置 4 0 に書き込み終了通知を返し（ステップ S 4 0 4）、ストレージ記憶部 2 0 6 がそのデータをストレージ装置ドライバ実行部 2 0 3 および外部ストレージ接続手段 2 3 を介して共有ストレージ 3 0 にそのデータを記憶する（ステップ S 4 0 5）。

30

## 【 0 0 6 1 】

共有ストレージ 3 0 への書き込みが終了したら、ミラーリング処理部 2 0 5 はそのデータを N V R A M 2 2 のグループ 1 用領域 2 2 a から削除する（ステップ S 4 0 6）。以上の処理によって、順序性をシビアに要求される処理であっても、データのリカバリ処理を行うことなく処理を継続することが可能となる。即ち、ごく短時間でサービスを復旧させることが可能となるので、クライアント装置 4 0 の側に故障の発生を意識させること自体が必要ない。

40

## 【 0 0 6 2 】

（実施形態の全体的な動作）

次に、上記の実施形態の全体的な動作について説明する。

本実施形態に係るフェイルオーバー方法は、同一の構成を有する他のノード装置と相互に接続されてクラスタシステムを構成すると共に、他のノード装置との間で共有される同一の外部記憶装置である共有ストレージに接続されてなるノード装置（第一ノード 1 0）にあって、予め備えられた不揮発性メモリは複数の記憶領域に区切られており、接続された各クライアント装置の I P アドレスと当該クライアント装置が使用すべき不揮発性メモ

50

りの記憶領域との間の対応関係が予め備えられたIPアドレステーブル記憶手段に記憶されたものであると共に、アプリケーションソフトをアプリケーションソフト実行部が実行し、アプリケーションソフトによって得られた処理データを不揮発性メモリ上のクライアント装置のIPアドレスに対応する記憶領域上にアプリケーションソフト実行部が一時的に保存し(図6・ステップS302または352)、記憶された処理データをミラーリング処理部が他のノード装置の対応する記憶領域に記憶させ(図6・ステップS303または353)、他のノード装置の対応する記憶領域に記憶させられた後の処理データをストレージ記憶部が共有ストレージに書き込む(図6・ステップS305または355)。

【0063】

そして、他のノード装置(第二ノード20)の不揮発性メモリに処理データが残った状態で当該他のノード装置に異常が発生した場合に、IPアドレステーブル記憶手段に記憶された対応関係をフェイルオーバー処理部が変更して当該他のノード装置による処理を引き継ぐ(図9・ステップS402~403)。

【0064】

ここで、上記各動作ステップについては、これをコンピュータで実行可能にプログラム化し、これらを前記各ステップを直接実行するノード装置(第一ノード10)のプロセッサ11に実行させるようにしてもよい。本プログラムは、非一時的な記録媒体、例えば、DVD、CD、フラッシュメモリ等に記録されてもよい。その場合、本プログラムは、記録媒体からコンピュータによって読み出され、実行される。

この動作により、本実施形態は以下のような効果を奏する。

【0065】

本実施形態によれば、第一ノード10もしくは第二ノード20のうちのいずれかに異常が発生して停止した場合でも、残る一方のノードが備えるNVRAM12または22に対応する記憶領域が存在しており、その領域に共有ストレージ30にすぐ書き込める状態のデータが、ミラーリング処理によって整合性が取れた状態で記憶されている。

【0066】

従って、IPアドレステーブル記憶手段15または25に記憶された各クライアント装置に属する記憶領域に応じて、データの書き込みを行う領域を切り替えれば、データのリカバリ処理を行う必要は無く、すぐに残る一方のノードで処理を引き継いで続行させることが可能となる。これによって、クライアント装置の側に故障の発生を意識させることなく、ごく短時間でサービスを復旧させることが可能となる。

【0067】

(実施形態の拡張)

上記実施形態は、以上で説明した本発明の趣旨を改変しない範囲で、様々な拡張が可能である。以下、これについて説明する。

【0068】

まず、上記実施形態は第一ノード10および第二ノード20という2台のサーバコンピュータによる構成例を示したが、これが3台以上になってももちろんよい。その場合、各クライアント装置をサーバコンピュータの台数分のグループに予め分けたIPアドレステーブルを、各々のIPアドレステーブル記憶手段に記憶させた上で、上記と同様の動作を行うこととなる。

【0069】

さらに、各サーバコンピュータが、物理的に複数台のコンピュータによって構成されてもよい。そして、各サーバコンピュータが、同一の処理をクライアント装置のグループ毎に分担してもよいし、また別々の処理を並行して行うものとしてもよい。

【0070】

これまで本発明について図面に示した特定の実施形態をもって説明してきたが、本発明は図面に示した実施形態に限定されるものではなく、本発明の効果を奏する限り、これまで知られたいかなる構成であっても採用することができる。

【0071】

10

20

30

40

50

上述した実施形態について、その新規な技術内容の要点をまとめると、以下のようになる。なお、上記実施形態の一部または全部は、新規な技術として以下のようにまとめられるが、本発明は必ずしもこれに限定されるものではない。

【 0 0 7 2 】

(付記 1) 同一の構成を有する他のノード装置と相互に接続されてクラスタシステムを構成すると共に、前記他のノード装置との間で共有される同一の外部記憶装置である共有ストレージに接続されてなるノード装置であって、

予め複数の記憶領域に区切られた不揮発性メモリと、

接続された各クライアント装置の IP アドレスと当該クライアント装置が使用すべき前記不揮発性メモリの前記記憶領域との間の対応関係を予め記憶している IP アドレステーブル記憶手段と、

前記クライアント装置からの要求に基づいて予め装備されたアプリケーションソフトを実行して処理を行い、これによって得られる処理データを処理依頼元の前記クライアント装置の IP アドレスに対応する前記記憶領域上に記憶させるアプリケーションソフト実行部と、

記憶された前記処理データを前記他のノード装置の対応する記憶領域に記憶させるミラーリング処理部と、

前記他のノード装置の対応する記憶領域に記憶させられた後の前記処理データを前記共有ストレージに書き込むストレージ記憶部と  
を有することを特徴とするノード装置。

【 0 0 7 3 】

(付記 2) 前記ミラーリング処理部が、前記処理データを前記他のノード装置の対応する記憶領域に記憶させる処理の完了後、処理要求元の前記クライアント装置に書き込み終了通知を返信する機能を有することを特徴とする、付記 1 に記載のノード装置。

【 0 0 7 4 】

(付記 3) 前記他のノード装置の不揮発性メモリに前記処理データが残った状態で当該他のノード装置に異常が発生した場合に、前記 IP アドレステーブル記憶手段に記憶された対応関係を変更して当該他のノード装置による処理を引き継ぐフェイルオーバー処理部を有することを特徴とする、付記 1 または付記 2 に記載のノード装置。

【 0 0 7 5 】

(付記 4) 同一の構成を有する第 1 および第 2 のノード装置と、前記第 1 および第 2 のノード装置の間で共有される同一の外部記憶装置である共有ストレージとが相互に接続されて構築されたクラスタシステムであって、

前記第 1 および第 2 のノード装置が、付記 1 ないし付記 3 のうちいずれか 1 項に記載のノード装置であることを特徴とするクラスタシステム。

【 0 0 7 6 】

(付記 5) 同一の構成を有する他のノード装置と相互に接続されてクラスタシステムを構成すると共に、前記他のノード装置との間で共有される同一の外部記憶装置である共有ストレージに接続されてなるノード装置にあって、

予め備えられた不揮発性メモリは複数の記憶領域に区切られており、

接続された各クライアント装置の IP アドレスと当該クライアント装置が使用すべき前記不揮発性メモリの前記記憶領域との間の対応関係が予め備えられた IP アドレステーブル記憶手段に記憶されたものであると共に、

前記アプリケーションソフトをアプリケーションソフト実行部が実行し、

前記アプリケーションソフトによって得られた前記処理データを前記不揮発性メモリ上の前記クライアント装置の IP アドレスに対応する前記記憶領域上に前記アプリケーションソフト実行部が一時的に保存し、

記憶された前記処理データをミラーリング処理部が前記他のノード装置の対応する記憶領域に記憶させ、

前記他のノード装置の対応する記憶領域に記憶させられた後の前記処理データをストレ

10

20

30

40

50

ージ記憶部が前記共有ストレージに書き込むことを特徴とするフェイルオーバー方法。

【 0 0 7 7 】

(付記6) 前記他のノード装置の不揮発性メモリに前記処理データが残った状態で当該他のノード装置に異常が発生した場合に、前記IPアドレステーブル記憶手段に記憶された対応関係をフェイルオーバー処理部が変更して当該他のノード装置による処理を引き継ぐ

ことを特徴とする、付記5に記載のフェイルオーバー方法。

【 0 0 7 8 】

(付記7) 同一の構成を有する他のノード装置と相互に接続されてクラスタシステムを構成すると共に、前記他のノード装置との間で共有される同一の外部記憶装置である共有ストレージに接続されてなるノード装置にあって、

予め備えられた不揮発性メモリは複数の記憶領域に区切られており、

接続された各クライアント装置のIPアドレスと当該クライアント装置が使用すべき前記不揮発性メモリの前記記憶領域との間の対応関係が予め備えられたIPアドレステーブル記憶手段に記憶されたものであると共に、

前記ノード装置が備えるプロセッサに、

前記アプリケーションソフトを実行する手順、

前記アプリケーションソフトによって得られた前記処理データを前記不揮発性メモリ上の前記クライアント装置のIPアドレスに対応する前記記憶領域上に一時的に保存する手順、

記憶された前記処理データを前記他のノード装置の対応する記憶領域に記憶させる手順、

および前記他のノード装置の対応する記憶領域に記憶させられた後の前記処理データを前記共有ストレージに書き込む手順

を実行させることを特徴とするフェイルオーバープログラム。

【産業上の利用可能性】

【 0 0 7 9 】

本発明は、コンピュータにおいてデータの整合性や順序性が要求される用途のクラスタシステムにおいて適用可能である。より具体的には、ファイルサーバ、データベースシステム、ウェブサーバ、業務システムなどに適用可能である。

【符号の説明】

【 0 0 8 0 】

- 1 クラスタシステム
- 10 第一ノード
- 11, 21 プロセッサ
- 12, 22 NVRAM
- 12a, 22a グループ0用領域
- 12b, 22b グループ1用領域1
- 13, 23 外部ストレージ接続手段
- 14, 24 通信手段
- 15, 25 IPアドレステーブル記憶手段
- 20 第二ノード
- 30 共有ストレージ
- 40 クライアント装置
- 41 ネットワーク
- 101, 201 アプリケーション実行部
- 102, 202 NVRAMドライバ実行部
- 103, 203 ストレージ装置ドライバ実行部
- 104, 204 外部通信ドライバ実行部

10

20

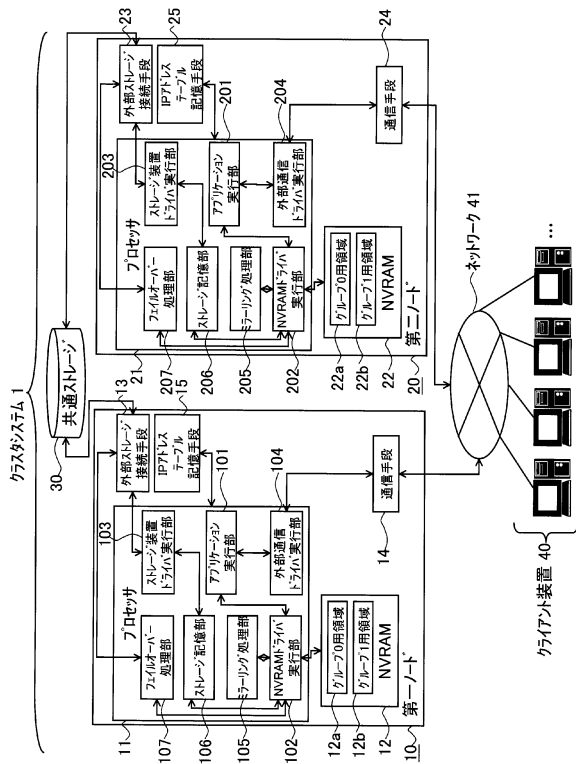
30

40

50

- 105, 205 ミラーリング処理部
- 106, 206 ストレージ記憶部
- 107, 207 フェイルオーバー処理部

【図1】

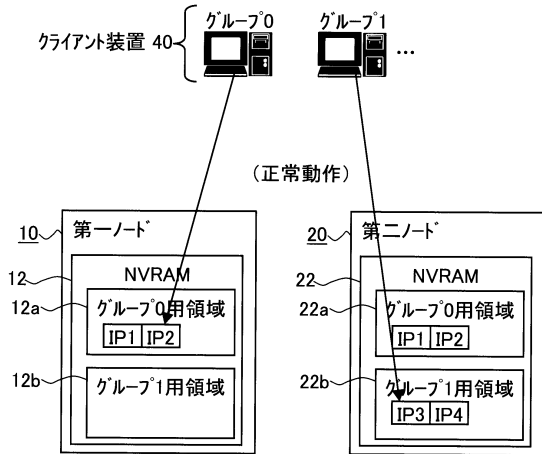


【図2】

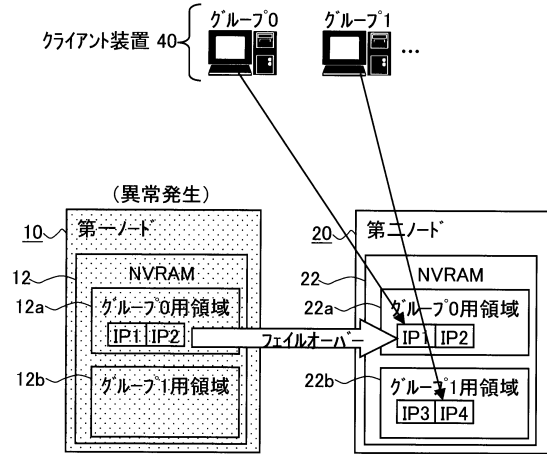
15 (25)

IPアドレス	グループ
192.168.10.1	0
192.168.10.2	1
192.168.10.10	1
...	...

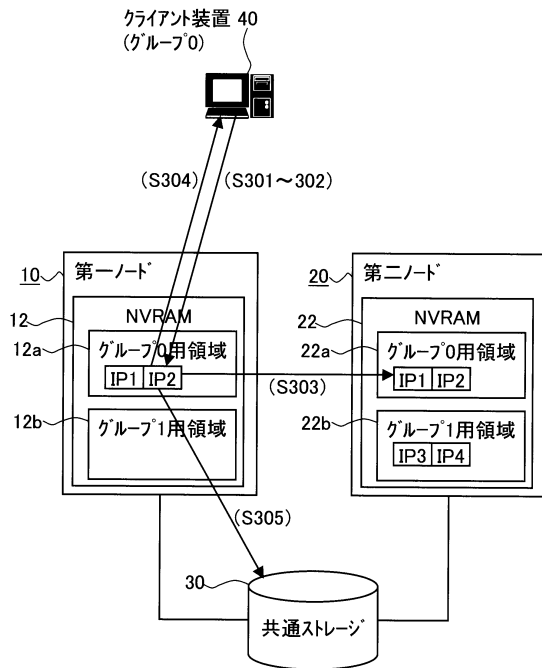
【図3】



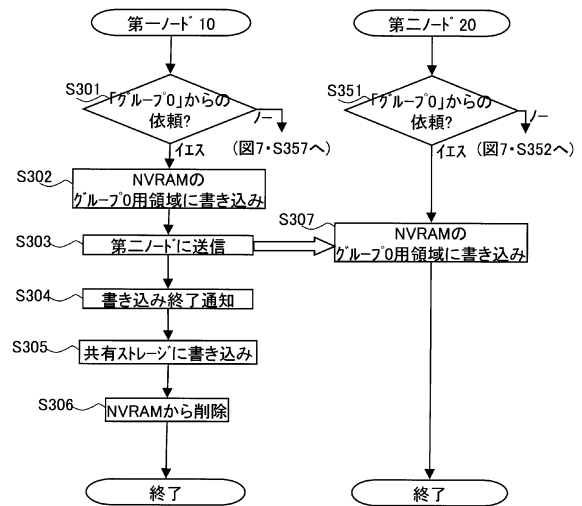
【図4】



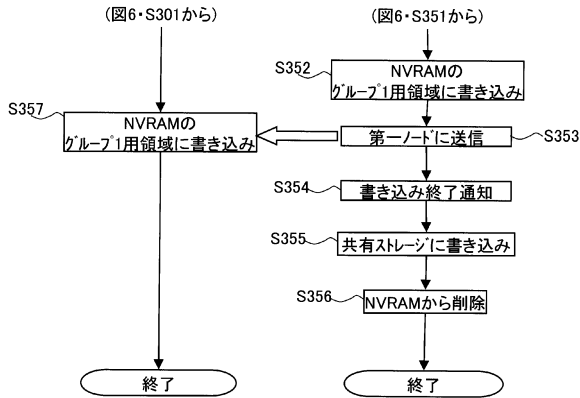
【図5】



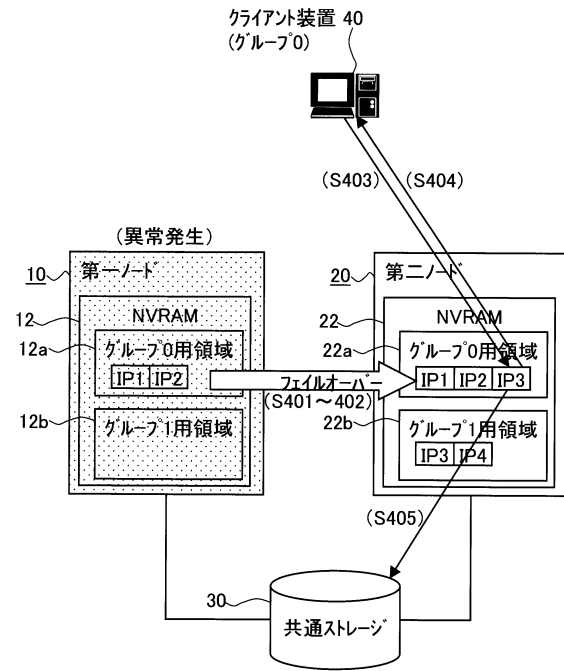
【図6】



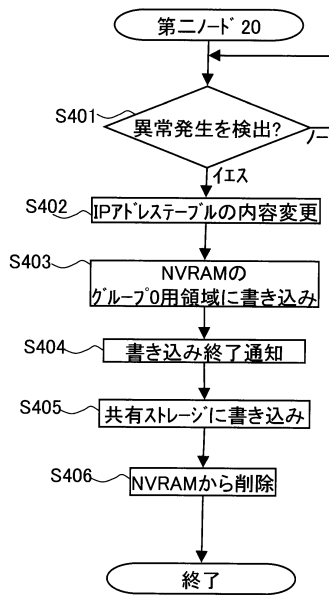
【図7】



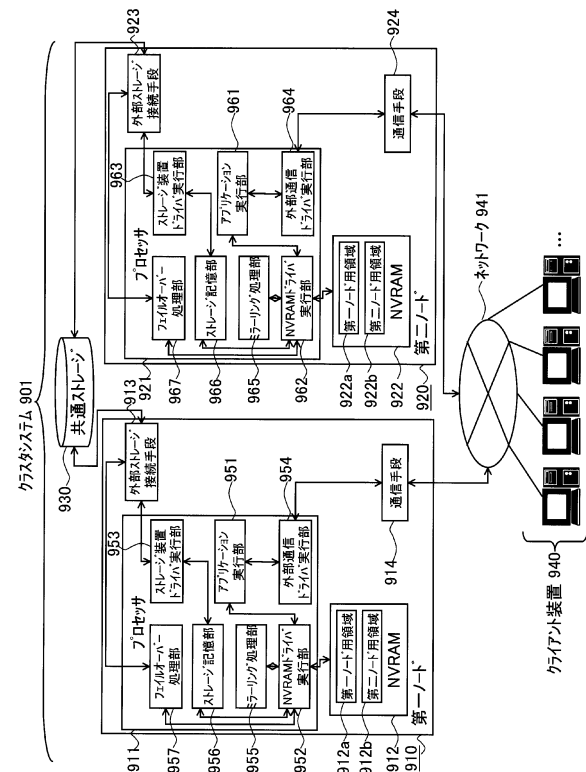
【図8】



【図9】



【図10】



---

フロントページの続き

- (56)参考文献 米国特許第07730153(US, B1)  
特開2002-373102(JP, A)  
特開2006-302153(JP, A)  
特開2001-175597(JP, A)  
国際公開第2006/057040(WO, A1)

- (58)調査した分野(Int.Cl., DB名)  
G06F 11/20  
G06F 9/50