

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5217966号  
(P5217966)

(45) 発行日 平成25年6月19日(2013.6.19)

(24) 登録日 平成25年3月15日(2013.3.15)

(51) Int.Cl. F 1  
**G 0 6 F 3/06 (2006.01)**  
 G 0 6 F 3/06 3 0 1 J  
 G 0 6 F 3/06 3 0 2 A

請求項の数 7 (全 23 頁)

(21) 出願番号	特願2008-304197 (P2008-304197)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成20年11月28日(2008.11.28)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2010-128885 (P2010-128885A)	(74) 代理人	100092152 弁理士 服部 毅巖
(43) 公開日	平成22年6月10日(2010.6.10)	(72) 発明者	野口 泰生 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
審査請求日	平成23年8月8日(2011.8.8)	(72) 発明者	内田 考介 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	坂東 博司

最終頁に続く

(54) 【発明の名称】 ストレージシステムのアップデート処理プログラム、アップデート処理方法及びストレージシステム

(57) 【特許請求の範囲】

【請求項1】

複数のストレージ装置に分散してデータを格納するストレージシステムにてコンピュータにアップデート処理を行わせるアップデート処理プログラムにおいて、

前記コンピュータを、

前記ストレージ装置へのデータ書き込み要求を受け取ったときに、書き込みデータを前記ストレージ装置に直接書き込む同期モードと、前記書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングで前記ストレージ装置に書き込む非同期モードとを備えた読み書き制御手段に対し、前記同期モードまたは前記非同期モードのどちらでデータ書き込みを行うかを指示する同期/非同期指示手段、

前記非同期モードで所定のサービス処理を実行中の第1のプロセスに代わって前記サービス処理を実行可能な第2のプロセスをサービス停止状態にして起動し、前記第1のプロセスから前記第2のプロセスへのアップデートが要求されると、前記第1のプロセスの前記サービス処理の実行を、前記同期/非同期指示手段に指示して前記非同期モードから前記同期モードに切り替え、前記第1のプロセスが前記同期モードで実行中の前記サービス処理を終了させ、前記第1のプロセスが前記サービス処理を終了した後に、前記第2のプロセスのサービス処理を開始させる、プロセス制御手段、

として機能させることを特徴とするストレージシステムのアップデート処理プログラム

【請求項2】

前記ストレージ装置の記憶領域は、仮想的な論理ボリュームを所定の記憶領域単位で分割した論理セグメントと、前記論理セグメントに対応付けられた前記ストレージ装置の実データ記憶領域を前記所定の記憶領域単位で分割したスライスとを関連付けた管理情報によって管理されており、

前記プロセス制御手段は、前記同期モードを指示する前に、前記複数のストレージ装置に分散して格納される前記管理情報を読み出し、前記第2のプロセスを、前記サービス処理を停止させた状態で起動しておく、

ことを特徴とする請求項1記載のストレージシステムのアップデート処理プログラム。

【請求項3】

前記コンピュータを、少なくとも、前記アップデート処理の開始前と、前記アップデート処理の終了時とに、前記ストレージ装置に分散して格納される、仮想的な論理ボリュームを所定の記憶領域単位で分割した論理セグメントと、前記論理セグメントに対応付けられた前記ストレージ装置の実データ記憶領域を前記所定の記憶領域単位で分割したスライスとを関連付けた管理情報を収集し、前記アップデート処理の開始前の前記管理情報と、前記アップデート処理の終了時の前記管理情報とを比較し、前記管理情報に変更があったか否かを前記プロセス制御手段に通知する管理情報チェック手段として機能させ、

前記プロセス制御手段は、前記管理情報の変更があったときは前記アップデート処理を中止し、前記アップデート処理が開始される前の状態に戻す、

ことを特徴とする請求項1記載のストレージシステムのアップデート処理プログラム。

【請求項4】

前記プロセス制御手段は、前記複数のストレージ装置それぞれに接続される複数の前記コンピュータを管理する管理装置から前記コンピュータに向けて一斉に送信される指示に応じて、前記読み書き制御手段を前記同期モードで動作させる処理、前記第1のプロセスに終了指示を出して処理を終了させる処理、及び前記第2のプロセスによるサービスを開始させる処理をそれぞれ実行し、処理が終了するごとに前記管理装置に終了結果を通知し、次の指示を待つ、

ことを特徴とする請求項1記載のストレージシステムのアップデート処理プログラム。

【請求項5】

前記プロセス制御手段は、1または複数の前記コンピュータが前記指示に応じた処理に失敗したときに前記管理装置によって発行されるロールバック指示を受けたときは、前記管理装置からの指示に応じて実行してきた前記アップデート処理の手順を逆順に辿り、前記アップデート処理が開始される前の状態に戻す、

ことを特徴とする請求項4記載のストレージシステムのアップデート処理プログラム。

【請求項6】

複数のストレージ装置に分散してデータを格納するストレージシステムにて、コンピュータがアップデート処理を行うアップデート処理方法において、

前記コンピュータが、

前記ストレージ装置へのデータ書き込み要求を受け取ったときに、書き込みデータを前記ストレージ装置に直接書き込む同期モードと、前記書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングで前記ストレージ装置に書き込む非同期モードとを備えた読み書き制御手段に対し、同期/非同期指示手段により前記同期モードまたは前記非同期モードのどちらでデータ書き込みを行うかを指示可能であって、

前記非同期モードで所定のサービス処理を実行中の第1のプロセスに代わって前記サービス処理を実行可能な第2のプロセスをサービス停止状態にして起動し、前記第1のプロセスから前記第2のプロセスへのアップデートが要求されると、前記第1のプロセスの前記サービス処理の実行を、前記同期/非同期指示手段に指示して前記非同期モードから前記同期モードに切り替え、前記第1のプロセスが前記同期モードで実行中の前記サービス処理を終了させ、前記第1のプロセスが前記サービス処理を終了した後に、前記第2のプロセスのサービス処理を開始させる、

アップデート処理を行うことを特徴とするストレージシステムのアップデート処理方法

10

20

30

40

50

。

【請求項 7】

ストレージ装置と接続してアクセス管理を行うストレージノードと、アップデート処理を管理する管理ノードとを備え、複数の前記ストレージ装置に分散してデータを格納するストレージシステムにおいて、

前記ストレージノードは、

前記ストレージ装置へのデータ書き込み要求を受け取ったときに、書き込みデータを前記ストレージ装置に直接書き込む同期モードと、前記書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングで前記ストレージ装置に書き込む非同期モードとを備えた読み書き制御手段に対し、前記同期モードまたは前記非同期モードのどちらでデータ書き込みを行うかを指示する同期 / 非同期指示手段と、

10

前記非同期モードで所定のサービス処理を実行中の第 1 のプロセスに代わって前記サービス処理を実行可能な第 2 のプロセスをサービス停止状態にして起動する起動処理を行い、前記第 1 のプロセスから前記第 2 のプロセスへのアップデートが要求されると、前記第 1 のプロセスの前記サービス処理の実行を、前記同期 / 非同期指示手段に指示して前記非同期モードから前記同期モードに切り替える切替処理を行い、前記第 1 のプロセスが前記同期モードで実行中の前記サービス処理を終了させる終了処理を行い、前記第 1 のプロセスが前記サービス処理を終了した後に、前記第 2 のプロセスのサービス処理を開始させる開始処理を行い、前記起動処理と前記切替処理と前記終了処理と前記開始処理が終了するごとに終了結果を前記管理ノードに通知する通知処理を行う、プロセス制御手段と、

20

を備え、

前記管理ノードは、

ネットワークを介して複数の前記ストレージノードに接続し、前記第 1 のプロセスを有する前記ストレージノードのアップデートの進行状況に関するアップデート管理情報を格納する管理情報記憶手段と、

前記管理情報記憶手段に格納される前記アップデート管理情報に基づき、対象ストレージノードすべてに対し前記同期モードへの切り替えを指示し、前記対象ストレージノードすべての処理が成功したときは前記対象ストレージノードに対し前記プロセスの停止を指示し、前記対象ストレージノードすべての処理が成功したときは前記対象ストレージノードに対し前記新プロセスを起動させ前記第 2 のプロセスによるサービスを開始させるアップデート管理手段と、

30

を備えることを特徴とするストレージシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はストレージシステムのアップデート処理プログラム、アップデート処理方法及びストレージシステムに関し、特に複数のストレージ装置に分散してデータを格納するストレージシステムにてコンピュータにアップデート処理を行わせるアップデート処理プログラム、アップデート処理方法及びそのストレージシステムに関する。

【背景技術】

40

【0002】

従来、ストレージシステムとして、複数のストレージノードをネットワーク上に分散配置して協働させることによって、性能及び信頼性を向上させる分散型のマルチノードストレージがある。マルチノードストレージでは、制御ノードが各ストレージノードを管理している。このため、各ストレージノードは、通信機能部を介して自身が動作していることを示す生存信号を定期的送信している。この生存信号をハートビートと呼ぶ。制御ノードはストレージノードから送出されるハートビートを監視しており、ストレージノードからのハートビートが途絶えた時は、リカバリ処理を行う（例えば、特許文献 1 参照）。

【0003】

マルチノードストレージでは、メンテナンスなどのため、運用中、システムを停止させ

50

ることなくストレージノードに搭載されるプロセスをアップデートする必要がある。図 1 4 は、従来のマルチノードストレージのプロセスアップデートの手順を示した図である。

【 0 0 0 4 】

図 1 4 の例のマルチノードストレージは、リカバリ処理などを行う制御ノード 9 0 と、データを分散して格納するディスクノード 9 1 , 9 2 , 9 3 と、を有する。ディスクノード 9 1 , 9 2 , 9 3 は、それぞれ一定周期でハートビート ( 図では H B ) を送出している。制御ノード 9 0 は、ディスクノード 9 1 , 9 2 , 9 3 からハートビートを受信し、各ディスクノード 9 1 , 9 2 , 9 3 の状態を検知する。

【 0 0 0 5 】

ここで、運用中にディスクノード 9 3 のアップデートを行う従来の手順について説明する。プロセスをアップデートするためには、ディスクノードの再起動が必要になる。そこで、ディスクノード 9 3 に対して旧プロセスの終了を指示し、旧プロセス終了処理 9 3 1 が行われる。続いて新プロセスで再起動する指示が出され、新プロセス再起動処理 9 3 2 が行われる。この処理が行われている間、ディスクノード 9 3 はハートビートを送出することができない。また、他装置からディスクノード 9 3 へのアクセスもできなくなる。そして、新プロセスの再起動処理 9 3 2 が終了後、ハートビートの送出手が再開され、通常状態に戻る。

【 0 0 0 6 】

一方、ディスクノード 9 3 のアップデートの間、制御ノード 9 0 は、ディスクノード 9 1 及びディスクノード 9 2 からハートビートを受信することができる。しかし、アップデート処理が開始されているディスクノード 9 3 からハートビートを受信できない。このディスクノード 9 3 の H B 途絶期間が一定時間を超過すると、制御ノード 9 0 は、ディスクノード 9 3 を故障と見なし、リカバリ処理 9 0 1 を行う。

【 特許文献 1 】再表 2 0 0 4 / 1 0 4 8 4 5 号公報

【 発明の開示 】

【 発明が解決しようとする課題 】

【 0 0 0 7 】

しかし、従来のマルチノードストレージには、運用中におけるプロセスのアップデート作業が容易ではないという問題点があった。

図 1 4 に示したように、プロセスをアップデートするためには、ディスクノードの再起動が必要となり、その間はハートビート ( H B ) が途絶する。また、この間は他装置からのアクセスに対しても応答することができない。さらに、そのハートビート途絶期間が一定時間を超過すると故障と判断され、リカバリが発生するという問題点もあった。このため、アップデートに要する時間を極力短縮させる必要がある。

【 0 0 0 8 】

しかし、新プロセスを稼働させるまでには、旧プロセス終了処理 9 3 1 と、新プロセス再起動処理 9 3 2 とを実行しなければならない。このうち、旧プロセス終了処理 9 3 1 では、O S ( Operating System、オペレーティングシステム ) がキャッシュメモリに一時保存しているデータをディスクに書き込んで同期させる処理が行われる。この同期処理は、キャッシュメモリに残っているデータの量に応じて処理時間が延び、キャッシュメモリに大量のデータが残っていた場合には非常に時間がかかることがある。また、新プロセス再起動処理 9 3 2 では、ディスクノードに分散配置されるデータに関するメタデータを読み込む必要がある。メタデータも他のディスクノードに分散配置されており、読み込むために時間がかかる。このように、旧プロセス終了処理にも新プロセス再起動処理にも処理時間を長期化する要素があり、アップデートの時間を短縮することは容易ではなかった。また、結果としてアップデートを行っているディスクノードが故障と見なされてしまうことを抑制することができなかった。

【 0 0 0 9 】

さらに、同時に 2 台以上のディスクノードが故障と判断されるときは、マルチノードス

10

20

30

40

50

トレージはシャットダウンしてしまう。このため、2台以上のディスクノードを同時にアップデートすることはできず、逐次的にアップデートしなければならないという問題点があった。このように1台ずつアップデートがされるため、全システムがアップデートされるまでには多大な時間がかかった。

【0010】

本発明はこのような点に鑑みてなされたものであり、運用中のプロセスのアップデートに要する時間を短縮させることが可能なマルチノードストレージシステムのアップデート処理プログラム、アップデート処理方法及びストレージシステムを提供することを目的とする。

【課題を解決するための手段】

【0011】

上記課題を解決するために、複数のストレージ装置に分散してデータを格納するストレージシステムにてコンピュータにアップデート処理を行わせるアップデート処理プログラムが提供される。このようなアップデート処理プログラムは、コンピュータを、同期/非同期指示手段と、プロセス制御手段と、して機能させる。同期/非同期指示手段は、ストレージ装置へのデータ書き込み要求を受け取ったときに、書き込みデータをストレージ装置に直接書き込む同期モードと、書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングでストレージ装置に書き込む非同期モードとを備えた読み書き制御手段に対し、同期モードまたは非同期モードのどちらでデータ書き込みを行うかを指示する。プロセス制御手段は、非同期モードで所定のサービス処理を実行中の第1のプロセスに代わってサービス処理を実行可能な第2のプロセスをサービス停止状態にして起動する。第1のプロセスから第2のプロセスへのアップデートが要求されると、第1のプロセスのサービス処理の実行を、同期/非同期指示手段に指示して非同期モードから同期モードに切り替える。こうして第1のプロセスが同期モードで実行中のサービス処理を終了させる。そして、第1のプロセスがサービス処理を終了した後に、第2のプロセスによるサービス処理を開始させる。

【0012】

このようなアップデート処理プログラムを実行するコンピュータは、第1のプロセスを第2のプロセスにアップデートする際に、サービス処理を実行中は非同期モードに設定される読み書き制御手段を同期モードにする。こうして第1のプロセスからの書き込みデータがストレージ装置に直接書き込まれている状態(同期モード)で第1のプロセスを終了させる。そして、第1のプロセスが処理を終了した後、第2のプロセスによるサービスを開始させる。

【0013】

また、上記課題を解決するために、上記のアップデート処理プログラムと同様の処理手順を実行させるアップデート処理方法、及び同様の機能を備えた装置を有するストレージシステムが提供される。

【発明の効果】

【0014】

開示のストレージシステムのアップデート処理プログラム、アップデート処理方法及びストレージシステムによれば、旧プロセス終了処理の前にストレージ装置へのデータ書き込みが同期モードで行われるようになる。これにより、プロセス終了処理において、キャッシュメモリに蓄積されているデータをディスクに書き込んで同期させる処理を省くことができ、アップデートに要する時間を短縮することができる。

【発明を実施するための最良の形態】

【0015】

以下、本発明の実施の形態を図面を参照して説明する。まず、発明の概要について説明し、その後、具体的な内容を説明する。

図1は、発明の概要を示した図である。図は、ストレージシステムを構成する複数のストレージノードのうちの1つを示している。

10

20

30

40

50

## 【 0 0 1 6 】

ストレージノード 1 0 は、データを格納するストレージ 2 0 に接続し、ネットワークを介して入力されるアクセスノードあるいは制御ノードから要求に応じてストレージ 2 0 内のデータへのアクセス管理を行う。

## 【 0 0 1 7 】

このストレージシステムは、論理ボリュームと呼ばれる仮想的なディスクを有する。この論理ボリュームを所定のサイズの単位（以下、セグメントとする）に分割し、各セグメントのデータを複数のストレージノードに分散配置する。ストレージ 2 0 のデータ記憶領域は、スライスと呼ばれる所定の単位に分割され、このスライスにセグメントが割り当てられる。この仮想的な論理ボリュームを分割した論理セグメントの識別情報と、この論理セグメントに対応付けられたストレージ装置の実データ記憶領域をセグメント単位で分割したスライスとを関連付けた管理情報は、メタデータと呼ばれる。具体的には、スライスに割り当てられたセグメントの情報（このセグメントの論理ディスク上のアドレスなど）などが記述され、セグメントデータとともに管理されている。

10

## 【 0 0 1 8 】

ストレージノード 1 0 は、ネットワークに接続する通信手段 1 1、運用中のプロセス 1 2、アップデート後の新プロセス 1 3、ストレージ 2 0 への読み書きを制御する読み書き制御手段（以下、R / W 制御手段とする）1 4、プロセスの動作を制御するプロセス制御手段 1 5、メタデータの変更をチェックするメタデータチェック手段 1 6 及びストレージ 2 0 への書き込みの動作モードを指示する同期 / 非同期指示手段 1 7 を有する。各処理手段は、コンピュータがそれぞれの処理を記述したプログラムを実行することにより、その機能を実現する。特に、プロセス制御手段 1 5、メタデータチェック手段 1 6 及び同期 / 非同期指示手段 1 7 は、コンピュータがアップデート処理プログラムを実行することにより、その処理機能を実現する。

20

## 【 0 0 1 9 】

通信手段 1 1 は、ストレージノード 1 0 と、図示しないネットワークを介して接続する制御ノード、アクセスノード及び他のストレージノードと、の間の通信を制御する。

プロセス 1 2 は、現在運用中のプロセスで、予め決められた所定のサービス処理を実行している。

## 【 0 0 2 0 】

新プロセス 1 3 は、アップデート後に、プロセス 1 2 の代わりにサービスを提供するプロセスであり、プロセス 1 2 と同様の機能を有する。

30

R / W 制御手段 1 4 は、プロセス 1 2 及び新プロセス 1 3 からのデータ読み出し要求またはデータ書き込み要求を受けて、ストレージ 2 0 からのデータ読み出し、またはストレージ 2 0 へのデータ書き込み処理を制御する。このうち、データの書き込み処理には、同期モードと非同期モードとがある。同期モードでは、データ書き込み要求を受け取ったとき、直接ストレージ 2 0 に書き込みを行い、書き込み終了の応答を返す。これに対し、非同期モードは、データ書き込み要求を受け取ったときに、キャッシュメモリにそのデータを蓄積し、書き込み終了の応答を返す。そして、所定のタイミングでキャッシュメモリに蓄積されたデータをストレージ 2 0 に書き込む。通常のプロセス処理では、応答性能を上げるため、非同期モードで書き込み処理が行われる。

40

## 【 0 0 2 1 】

プロセス制御手段 1 5 は、通信手段 1 1 を介して入力される管理ノードの指示などに基づいて、管理下のプロセスの動作を制御する。アップデートの際には、サービスを停止した状態で新プロセス 1 3 を起動させる。次に、プロセス 1 2 がサービス処理を行っている状態で、同期 / 非同期指示手段 1 7 に指示し、R / W 制御手段 1 4 を同期モードに設定する。同期モードで動作している状態で、プロセス 1 2 の終了処理を指示し、プロセス 1 2 を終了させる。その後、新プロセス 1 3 のサービスを開始させる。このとき、同期 / 非同期指示手段 1 7 に指示し、R / W 制御手段 1 4 を非同期モードに戻す。

## 【 0 0 2 2 】

50

メタデータチェック手段16は、プロセス制御手段15からの指示に応じて自ノードのメタデータ変更チェックを行う。少なくとも、アップデート処理の開始前と、アップデート処理終了時のメタデータを比較する。こうして、アップデート処理の前後でメタデータに変更があったか否かを判定し、判定結果をプロセス制御手段15に通知する。

【0023】

同期/非同期指示手段17は、プロセス制御手段15の指示に従って、R/W制御手段14に対して同期モード/非同期モードの切り替え指示を行う。

また、プロセスが正常に動作しているときは、図示しないハートビート送出手段によって、他装置に向けてハートビートが送出される。ハートビートは、自装置の状態を示す生存信号であり、定期的に送出される。

【0024】

このような構成のストレージノード10によるプロセスのアップデート動作について説明する。現在、プロセス12がサービス処理を実行している。R/W制御手段14は、非同期モードで動作している。非同期モードとすることにより、書き込み時の応答時間を短縮することができる。

【0025】

アップデートする新プロセス13のプログラムは、事前に記憶手段にロードしておく。プロセス制御手段15は、メタデータチェック手段16を介して他装置のメタデータを読み出し、サービスを停止した状態で新プロセス13を起動する。新プロセス13は、サービス処理は停止しているが、メタデータを読み込んだ後の起動処理を行う。なお、メタデータチェック手段16は、このときのメタデータを記憶しておく。次に、同期/非同期指示手段17に指示し、R/W制御手段14を同期モードにする。これにより、プロセス12がデータの書き込み要求を行うと、すぐにストレージ20に書き込まれるようになる。続いて、プロセス制御手段15は、プロセス12に対して処理の終了を指示する。R/W制御手段14は、同期モードで動作しているので、キャッシュメモリのデータをストレージ20に書き込む同期処理は必要ない。このため、プロセス12は、終了処理を短時間で完了することができる。プロセス12を終了させた後、メタデータチェック手段16により、メタデータが変更されていないかどうかをチェックする。メタデータチェック手段16は、再度メタデータを収集し、新プロセス13を起動したときのメタデータと比較する。なお、比較は、処理ごとに行われてもよい。メタデータが変更されていないときは、プロセス制御手段15は、新プロセス13を再起動する。新プロセス13は、起動時のメタデータ読み込みは終了しているので、再起動処理を短時間で完了し、サービス処理を開始する。一方、メタデータが変更されていた場合には、アップデートを中断し、一旦、アップデート開始前の状態(非同期モードでプロセス12がサービス処理を行っている状態)に戻す。そして、必要であれば、新たなメタデータで、再度上記手順でアップデート処理を行う。

【0026】

このように、ストレージノード10によれば、プロセス12から新プロセス13へのアップデート時、プロセス12の終了処理時間を延ばしていたデータの同期化処理を行う必要がない。また、新プロセス13の再起動に要する時間を延ばしていたメタデータの読み出し処理を行う必要がない。これにより、アップデートに要する時間を短縮することができる。この結果、再起動時、制御ノードが故障と判定する前にハートビートの送出を開始することが可能となり、リカバリ処理を抑止することができる。

【0027】

以下、発明の実施の形態を詳細に説明する。図2は、本実施の形態のマルチノードストレージの構成例を示す図である。

マルチノードストレージは、ネットワーク400を介して、複数のディスクノード100, 200, 300と、アクセスノード500、制御ノード600及び管理ノード700が接続されている。

【0028】

ディスクノード100にはディスク110、ディスクノード200にはディスク210、ディスクノード300にはディスク310が、それぞれ接続されている。ディスク110には、複数のハードディスク装置(HDD)が実装されている。ディスク210、310の構成も同様である。ディスクノード100、200、300は、例えば、IA(Intel Architecture)と呼ばれるアーキテクチャのコンピュータである。そして、接続されたディスク110、210、310に格納されたデータを管理し、管理しているデータをアクセスノード500経由で端末装置801、802、803に提供する。また、ディスクノード100、200、300は、冗長性を有するデータを管理することもできる。この場合、同一のデータが、少なくとも2つのディスクノードで管理される。本実施の形態では、ディスクノード100、200、300として、図1に示したアップデート処理を行うストレージノードを提供する。

10

**【0029】**

アクセスノード500には、ネットワーク800を介して複数の端末装置801、802、803が接続されている。アクセスノード500は、ディスクノード100、200、300のそれぞれが管理しているデータの格納場所を認識しており、端末装置801、802、803からの要求に回答して、ディスクノード100、200、300へデータアクセスを行う。

**【0030】**

制御ノード600は、ディスクノード100、200、300を管理する。例えば、制御ノード600は、ディスクノード100、200、300から送出されるハートビートを監視し、故障を検出したときはリカバリ処理を行う。

20

**【0031】**

管理ノード700は、マルチノードストレージのシステム全体を管理する。例えば、管理者からのアップデート指示に応じて、ディスクノード100、200、300全体のアップデート処理を管理する。

**【0032】**

図3は、ディスクノードのハードウェア構成例を示す図である。ディスクノード100は、CPU(Central Processing Unit)101によって装置全体が制御されている。CPU101には、バス106を介してRAM(Random Access Memory)102、HDD103、通信インタフェース104及びHDDインタフェース105が接続されている。

30

**【0033】**

RAM102には、CPU101に実行させるOSやアプリケーションプログラムの少なくとも一部が一時的に格納される。また、RAM102には、CPU101による処理に必要な各種データが格納される。HDD103には、OSやアプリケーションのプログラムが格納される。通信インタフェース104は、ネットワーク400に接続されている。通信インタフェース104は、ネットワーク400を介して、他のディスクノード、アクセスノード500、制御ノード600及び管理ノード700など、マルチノードストレージを構成する他のコンピュータとの間でデータの送受信を行う。HDDインタフェース105は、ディスク110を構成するHDDへのアクセス処理を行う。

**【0034】**

以上のようなハードウェア構成によって、本実施の形態の処理機能を実現することができる。なお、図3には、ディスクノード100を示したが、他のディスクノード200、300も同様のハードウェア構成で実現される。

40

**【0035】**

次に、上記のマルチノードストレージにおいてアップデート処理を行う各部について説明する。図4は、マルチノードストレージにおいてアップデート処理を行う各部のソフトウェア構成を示した図である。

**【0036】**

ディスクノード100、200、300は、スイッチ(SW)401を介して、アクセスノード500、制御ノード600及び管理ノード700との間でデータ交換を行う。こ

50



のうち、ディスクノード100, 200, 300のアップデート管理は管理ノード700が行う。

【0037】

ディスクノード100は、プロセス112及び新プロセス113と、アップデート処理を行うエージェント115とを有する。また、DP1というIDが付与されている。同様に、ディスクノード200は、プロセス212及び新プロセス213と、エージェント215とを有し、DP2というIDが付与されている。ディスクノード300の構成も同様であり、DP3というIDが付与されている。

【0038】

プロセス112, 212は、アップデート前のプロセスで現在サービス処理を行っているプロセスである。新プロセス113, 213は、アップデート後のプロセスである。プロセス112, 212及び新プロセス113, 213は、エージェント115, 215の指示に従って動作し、処理を全く実行していない状態、サービス処理以外の処理を実行している状態及びサービス処理を実行している状態のいずれかの状態にある。例えば、プロセス起動指示で、処理を全く実行していない状態からサービス処理以外の処理を実行している状態に遷移する。サービス開始指示でサービス処理を実行している状態に遷移する。サービス停止指示でサービス処理を実行している状態からサービス処理以外の処理を実行している状態に遷移する。そして、終了指示ですべての処理を終了し、処理を全く実行していない状態に遷移する。

【0039】

エージェント115, 215は、管理ノード700からの指示に従って、プロセス112, 212を、新プロセス113, 213にアップデートする。このため、プロセス制御手段、メタデータチェック手段及び同期/非同期指示手段としての機能を有する。なお、ディスク110へのアクセスを同期モードで行うか、非同期モードで行うかは、OSが管理する場合が多い。通常OSには、モードを切り替えるためのコマンドが用意されており、エージェント115, 215は、OSに対しこのようなコマンドを出力し、同期/非同期モードの切り替えを行う。例えば、非同期モードから同期モードへの切り替えは、アクセス受付を一時停止、デバイスファイルをクローズ、デバイスファイルを同期モードでオープン、アクセス受付を再開、という手順でコマンドを出力して行う。また、同期モードから非同期モードへの切り替えも同様に、アクセス受付を一時停止、デバイスファイルをクローズ、デバイスファイルを非同期モードでオープン、アクセス受付を再開、という手順でコマンドを出力して行う。

【0040】

管理ノード700は、アップデートを管理するアップデート管理部701と、管理テーブル702とを有する。管理テーブル702は、管理下のディスクノード100, 200, 300のアップデート進行状況を管理するためのアップデート管理情報が設定される。例えば、発行コマンドとその結果とが、ディスクノードごとに管理される。アップデート管理部701は、管理テーブル702に基づいて、ディスクノード100のエージェント115、ディスクノード200のエージェント215及び図示しないディスクノード300のエージェントと通信を行って、管理下のディスクノードのアップデートを一斉に処理する。

【0041】

このため、アップデート管理部701は、管理対象のディスクノード100, 200, 300に対し、アップデートの手順に沿った指示をコマンドとして順次出力する。例えば、新プロセスの起動を指示するコマンドを管理下のディスクノード100, 200, 300に出力する。そして、ディスクノード100, 200, 300のすべてから正常終了の応答が得られたときは、書き込み制御の同期モードへの変更を指示するコマンドをディスクノード100, 200, 300に出力する。同様にして、プロセスの終了、新プロセスのサービス開始、などコマンドを順次出力する。ディスクノード100, 200, 300のエージェントが、コマンドを受けて処理を行うことにより、すべてのディスクノードの

10

20

30

40

50

アップデートが同時に行われる。

【 0 0 4 2 】

以下、管理ノード 7 0 0 による一斉アップデート処理について詳細に説明する。まず、管理テーブル 7 0 2 について説明する。

図 5 は、管理テーブルの一例を示した図である。

【 0 0 4 3 】

管理テーブル 7 0 2 には、管理対象のディスクノードの管理情報として、ノード I D 7 0 2 1、コマンド発行 7 0 2 2 及び結果 7 0 2 3 の情報項目が登録される。

ノード I D 7 0 2 1 には、管理対象のディスクノードの識別情報（ディスクノードに付与された I D）が登録される。ここでは、ディスクノード 1 0 0 の D P 1、ディスクノード 2 0 0 の D P 2、ディスクノード 3 0 0 の D P 3 が登録される。

10

【 0 0 4 4 】

コマンド発行 7 0 2 2 には、コマンドの発行状態が登録される。例えば、コマンドをディスクノード 1 0 0、2 0 0、3 0 0 に出力していないときは、無（N U L L）が登録される。そして、コマンドを発行したときには、コマンドを送信したディスクノードに対応する欄に、発行済（D O N E）が登録される。なお、発行したコマンドの種別を登録してもよい。

【 0 0 4 5 】

結果 7 0 2 3 には、コマンド発行後に、ディスクノード 1 0 0、2 0 0、3 0 0 から得られた応答に基づいて、ディスクノード 1 0 0、2 0 0、3 0 0 の処理結果が登録される。結果を受け取るまでは、無（N U L L）が登録される。応答を受け取り、その結果が正常終了であれば、完了（O K）が登録される。そして、応答を受け取り、その結果が正常終了でなかったときは、失敗（N G）が登録される。

20

【 0 0 4 6 】

図 6 は、コマンドの発行処理に応じた管理テーブルの変化を示した図である。（A）はコマンド発行前、（B）はコマンド発行後、（C）はコマンド発行後結果受け付け中、（D）はコマンド結果受信後、（E）はコマンド結果受信後（N G 含む）の状態を示している。

【 0 0 4 7 】

（A）コマンド発行前は、すべてのディスクノード（D P 1、D P 2、D P 3）について、コマンド発行 7 0 2 2 が無（N U L L）の状態になっている。

30

（B）コマンド発行後は、すべてのディスクノード（D P 1、D P 2、D P 3）について、コマンド発行 7 0 2 2 が発行済（D O N E）の状態になっている。このように、コマンドの発行は、管理対象のディスクノードに対して一斉に行われる。

【 0 0 4 8 】

（C）コマンド発行後結果受け付け中は、コマンドを発行したディスクノード（D P 1、D P 2、D P 3）からの処理結果の応答を待っている状態である。ディスクノード（D P 1、D P 2、D P 3）は、発行されたコマンドの処理が終了すると、その処理結果を管理ノード 7 0 0 に向けて送信する。この例では、ディスクノード D P 1 と、ディスクノード D P 3 からの応答が得られ、ディスクノード D P 2 からの応答が得られていない状態である。

40

【 0 0 4 9 】

（D）コマンド結果受信後は、すべてのディスクノード（D P 1、D P 2、D P 3）から応答が得られた状態である。この例は、ディスクノード D P 1、ディスクノード D P 2 及びディスクノード D P 3 のすべてから正常完了（O K）が得られたことを示している。

【 0 0 5 0 】

（E）コマンド結果受信後（N G 含む）は、すべてのディスクノード（D P 1、D P 2、D P 3）から応答が得られ、応答に失敗（N G）が含まれていた場合である。この例は、ディスクノード D P 1 と、ディスクノード D P 3 からは正常完了（O K）の応答が得られ、ディスクノード D P 2 からの失敗（N G）の応答が得られたことを示している。

50

## 【 0 0 5 1 】

管理ノード700は、この管理テーブル702に基づいてアップデート処理を進める。

次に、上記の構成のマルチノードストレージのアップデート処理動作及びアップデート方法について詳しく説明する。

## 【 0 0 5 2 】

図7は、アップデート処理の動作シーケンス（同期モードへの切り替えまでの手順）を示した図である。図7は、管理ノード700からの指示がすべて成功した場合の例である。また、エージェント*i*は、任意のディスクノードに搭載されるエージェントであり、プロセス*i*はその運用中のプロセス、新プロセス*i*はその新プロセスを表している。

## 【 0 0 5 3 】

システム管理者などによって、管理ノード700にアップデート指示が出されると、処理が開始される。なお、新プロセス*i*のプログラムは、事前に各ディスクノードにダウンロードされているとする。

## 【 0 0 5 4 】

管理ノード700から、管理下のディスクノードすべてに対し、新プロセス起動のコマンド（1001）が出力される。コマンドを受け取ったエージェント*i*は、新プロセス起動を新プロセス*i*に指示する。これにより、新プロセス*i*において起動処理（1002）が行われる。起動処理（1002）では、メタデータを読み込み、サービス処理は停止した状態で新プロセス*i*を起動する。起動処理が終了し、サービス処理以外の処理が実行状態となった後、起動処理正常完了（OK）がエージェント*i*に返る。失敗時には、失敗（NG）が返る。正常完了（OK）を取得したエージェント*i*は、メタデータ変更チェック（1003）を行う。メタデータ変更チェック（1003）では、メタデータを収集し、このメタデータと、新プロセス*i*が起動処理（1002）で用いたメタデータ（前回収集したメタデータ）とが同じであるかどうかを判定する。メタデータが変更されていなければ、正常完了（OK）の応答を管理ノード700に返す。メタデータが変更されていれば、失敗（NG）の応答を管理ノード700に返す。なお、以下のメタデータ変更チェックでも、同様のチェックが行われる。ここでは、正常完了（OK）が返ったとして説明を続ける。エージェント*i*から正常完了（OK）を受けとった管理ノード700は、他のエージェントからの応答を待ち、全ノードが正常に終了したかどうかをチェックする（1004）。全ノードが新プロセス起動を正常完了したとき、次の手順へ処理を進める。一部のノードで新プロセスの起動に失敗、またはメタデータの変更が発生したとき、ロールバック処理を行う。ロールバック処理については、後述する。

## 【 0 0 5 5 】

全ノードが新プロセス起動を正常に完了したときは、管理ノード700から同期モード指示（1011）が出される。コマンドを受け取ったエージェント*i*は、同期モードを指示する。これにより、現在動作中のプロセス*i*のディスクへの書き込みモードが、同期モードへ変更（1012）される。同期モードへの変更の正常完了（OK）がエージェント*i*に返る。失敗時には、失敗（NG）が返る。正常完了（OK）を取得したエージェント*i*は、メタデータ変更チェック（1013）を行う。メタデータが変更されていなければ、正常完了（OK）の応答を管理ノード700に返す。メタデータが変更されていれば、失敗（NG）の応答を管理ノード700に返す。ここでは、正常完了（OK）が返ったとして説明を続ける。エージェント*i*から正常完了（OK）を受けとった管理ノード700は、他のエージェントからの応答を待ち、全ノードが正常に終了したかどうかをチェックする（1014）。全ノードが新プロセス起動を正常完了したとき、次の手順へ処理を進める。一部のノードでプロセスの同期化に失敗、またはメタデータの変更が発生したとき、ロールバック処理を行う。

## 【 0 0 5 6 】

次に、図8を用いて説明する。図8は、アップデート処理の動作シーケンス（新プロセスへの切り替えまでの手順）を示した図である。

全ノードが同期化モードへの移行を正常完了したときは、管理ノード700からプロセ

10

20

30

40

50

ス i のサービス停止指示 ( 1 0 2 1 ) が出される。コマンドを受け取ったエージェント i は、プロセス i に対してサービス停止を指示する。これにより、現在動作中のプロセス i が、サービス停止処理 ( 1 0 2 2 ) を行い、サービスを停止させる。サービス停止とともに、ハートビートの送信も停止され、ハートビート ( H B ) 途絶期間が開始される。プロセス i のサービス停止が完了したときは、正常完了 ( O K ) がエージェント i に返る。失敗時には、失敗 ( N G ) が返る。正常完了 ( O K ) を取得したエージェント i は、メタデータ変更チェック ( 1 0 2 3 ) を行う。メタデータが変更されていなければ、正常完了 ( O K ) の応答を管理ノード 7 0 0 に返す。メタデータが変更されていれば、失敗 ( N G ) の応答を管理ノード 7 0 0 に返す。ここでは、正常完了 ( O K ) が返ったとして説明を続ける。エージェント i から正常完了 ( O K ) を受けとった管理ノード 7 0 0 は、他のエージェントからの応答を待ち、全ノードが正常に終了したかどうかをチェックする ( 1 0 2 4 ) 。全ノードがプロセス i のサービス停止を正常完了したとき、次の手順へ処理を進める。一部のノードでプロセスの同期化に失敗、またはメタデータの変更が発生したとき、ロールバック処理を行う。

#### 【 0 0 5 7 】

全ノードがプロセス i のサービス停止処理を正常完了したときは、管理ノード 7 0 0 から新プロセス i への切り替え指示 ( 1 0 3 1 ) が出される。コマンドを受け取ったエージェント i は、新プロセス i に対して切り替え指示 ( サービス開始指示 ) を行う。これにより、現在サービス処理を停止中の新プロセス i が、サービス開始処理 ( 1 0 3 2 ) を行い、サービスを開始する。新プロセス i は、起動処理 1 0 0 2 によって既に動作を開始しているため、すぐにサービス処理を介することができる。サービス開始とともに、ハートビートの送信も再開され、ハートビート ( H B ) 途絶期間が終了する。新プロセス i のサービス開始が完了したときは、正常完了 ( O K ) がエージェント i に返る。失敗時には、失敗 ( N G ) が返る。正常完了 ( O K ) を取得したエージェント i は、サービスを停止しているプロセス i に対し、終了指示を行い、プロセス i において終了処理 ( 1 0 3 3 ) が行われる。プロセス i の終了処理が正常に完了していれば、正常完了 ( O K ) の応答を管理ノード 7 0 0 に返す。

#### 【 0 0 5 8 】

このように、プロセス i のサービス停止処理 1 0 2 2 が開始されるときには、既に同期化が終了しているため、同期化に要する時間がなくなる。また、新プロセス i のサービス開始 1 0 3 2 においても、既にメタデータを用いた起動処理 1 0 0 2 は終了しているため、すぐにサービスを開始できる。この結果、プロセスのアップデートに要する時間を大幅に短縮することができる。

#### 【 0 0 5 9 】

次に、ロールバック処理について説明する。

図 9 は、新プロセス起動に失敗したときの動作シーケンスを示した図である。図 9 は、( a ) 一部のノードで新プロセスの起動に失敗またはメタデータ変更発生、が起きた場合のロールバック処理である。

#### 【 0 0 6 0 】

新プロセス起動処理後の全ノード終了チェック ( 1 0 0 4 ) において、1 またはそれ以上のエージェントから失敗 ( N G ) の応答を受け取ったことが検出されたときは、管理ノード 7 0 0 は、ロールバック指示 ( 1 0 4 1 ) を管理対象のすべてのエージェント i へ出力する。ロールバック処理は、一連の処理が開始される前の状態、ここでは、新プロセス起動処理の指示が出される前、すなわち、アップデート開始前の状態に戻す処理を言う。ロールバック指示 ( 1 0 4 1 ) を受け取ったエージェント i は、起動した新プロセスを終了させるため、プロセス終了指示 ( 1 0 4 2 ) を新プロセス i へ出力する。新プロセス i は、プロセス終了処理 ( 1 0 4 3 ) を実行し、すべての処理を停止させる。その後、プロセスが正常に終了 ( O K ) したことをエージェント i に通知し、エージェント i が正常終了 ( O K ) したことを通知する応答を管理ノード 7 0 0 に返す。

#### 【 0 0 6 1 】

このような処理が行われることにより、一部のディスクノードで新プロセスの起動に失敗、またはメタデータ変更が発生したときは、アップデートが開始される前の状態にすべてのディスクノードが戻される。

【 0 0 6 2 】

図 1 0 は、プロセスの同期化に失敗したときの動作シーケンスを示した図である。図 1 0 は、( b ) 一部のノードでプロセスの同期化に失敗またはメタデータ変更発生、が起きた場合のロールバック処理である。

【 0 0 6 3 】

同期モードへの変更処理後の全ノード終了チェック ( 1 0 1 4 ) において、1 またはそれ以上のエージェントから失敗 ( N G ) の応答を受け取ったことが検出されたときは、管理ノード 7 0 0 は、ロールバック指示 ( 1 0 5 1 ) を管理対象のすべてのエージェント  $i$  10  
に出力する。ロールバック指示 ( 1 0 5 1 ) を受け取ったエージェント  $i$  は、同期モードを非同期モードに戻すため、非同期指示 ( 1 0 5 2 ) を行う。これにより、現在動作中のプロセス  $i$  のディスクへの書き込みモードは、非同期モードに変更 ( 1 0 5 3 ) される。その後、書き込みモードが非同期モードになったこと ( O K ) がエージェント  $i$  に通知されたときは、プロセス終了指示 ( 1 0 5 4 ) を新プロセス  $i$  に出力する。新プロセス  $i$  は、プロセス終了処理 ( 1 0 5 5 ) を実行し、すべての処理を停止させる。その後、プロセスが正常に終了 ( O K ) したことをエージェント  $i$  に通知し、エージェント  $i$  が正常終了 ( O K ) したことを通知する応答を管理ノード 7 0 0 に返す。

【 0 0 6 4 】

このような処理が行われることにより、一部のノードでプロセスの同期化に失敗またはメタデータ変更が発生したときも、アップデートが開始される前の状態にすべてのディスクノードが戻される。

【 0 0 6 5 】

図 1 1 は、プロセスのサービス停止に失敗したときの動作シーケンスを示した図である。図 1 1 は、( c ) 一部のノードでプロセスのサービス停止に失敗、またはメタデータ変更発生、が起きた場合のロールバック処理である。

【 0 0 6 6 】

プロセス  $i$  のサービス停止処理後の全ノード終了チェック ( 1 0 2 4 ) において、1 またはそれ以上のエージェントから失敗 ( N G ) の応答を受け取ったことが検出されたときは、管理ノード 7 0 0 は、ロールバック指示 ( 1 0 6 1 ) を管理対象のすべてのエージェント  $i$  30  
に出力する。ロールバック指示 ( 1 0 6 1 ) を受け取ったエージェント  $i$  は、サービスを停止させたプロセス  $i$  にサービス再開指示 ( 1 0 6 2 ) を行う。プロセス  $i$  は、サービス再開処理 ( 1 0 6 3 ) を実行し、サービスを再開する。サービス再開が正常完了 ( O K ) したことがプロセス  $i$  より通知されると、エージェント  $i$  は、非同期指示 ( 1 0 6 4 ) を行い、書き込みモードを非同期モードにする。非同期モードに変更 ( 1 0 6 5 ) され、正常完了 ( O K ) を受け取ると、新プロセス  $i$  に対し、プロセス終了指示 ( 1 0 6 7 ) を行う。新プロセス  $i$  は、プロセス終了処理 ( 1 0 6 8 ) を実行し、すべての処理を停止させる。その後、プロセスが正常に終了 ( O K ) したことをエージェント  $i$  に通知し、エージェント  $i$  が正常終了 ( O K ) したことを通知する応答を管理ノード 7 0 0 に返す。 40

【 0 0 6 7 】

このような処理が行われることにより、一部のノードでプロセスのサービス停止に失敗またはメタデータ変更が発生したときも、アップデートが開始される前の状態にすべてのディスクノードが戻される。

【 0 0 6 8 】

このように、アップデートの途中でエラーやメタデータの変更が生じたときは、その時点までに実行した処理を逆順に辿り、アップデート前の状態に戻される。これにより、アップデート未完の場合、プロセスがプロセスに戻ってサービスが継続される。

【 0 0 6 9 】

次に、管理ノード及びディスクノード ( エージェント ) のアップデート処理の手順を、

10

20

30

40

50

フローチャートを用いて説明する。

図12は、管理ノードのアップデート処理手順を示したフローチャートである。

【0070】

システム管理者などからの一斉アップデート指示を受け、管理ノード700が処理を開始させる。

【ステップS01】 新プロセス起動指示のコマンドを管理対象のすべてのディスクノードに向けて送信する。新プロセス起動指示は、新プロセスを、サービスを停止させた状態で起動させるための指示である。指示を受けたエージェントは、メタデータを読み込んで新プロセスの起動処理を行い、すぐにサービス開始ができる状態になった後、管理ノード700に応答を返してくる。

10

【0071】

【ステップS02】 新プロセス起動指示を出力した全ディスクノードからの応答を待つ。得られた応答をチェックし、全ディスクノードで新プロセス起動が成功したかどうかを判定する。成功したときは、処理をステップS03に進める。一部のディスクノードが、新プロセスの起動に失敗、またはメタデータ変更が発生して、失敗(NG)応答を返してきたときは、処理をステップS09に進める。

【0072】

【ステップS03】 新プロセス起動処理が正常に終了したときは、続いて、ディスクへの書き込みモードを同期モードとするように全ディスクノードに対し指示を出す。指示を受けたディスクノードは、書き込みモードを同期モードに切り替え、管理ノード700

20

【0073】

【ステップS04】 新プロセス起動指示を出力した全ディスクノードからの応答を待つ。得られた応答をチェックし、全ディスクノードで同期モードへの切り替えが成功したかどうかを判定する。成功したときは、処理をステップS05に進める。一部のディスクノードが、プロセスの同期化に失敗、またはメタデータ変更が発生して、失敗(NG)応答を返してきたときは、処理をステップS09に進める。

【0074】

【ステップS05】 同期モードへの変更が正常に終了したときは、プロセスのサービスを停止させる指示を出す。指示を受けたディスクノードは、プロセスに対し、プロセス終了を指示する。プロセスがプロセスを終了したら、管理ノード700に応答を返してくる。

30

【0075】

【ステップS06】 プロセスへのサービス停止指示を出力した全ディスクノードからの応答を待つ。得られた応答をチェックし、全ディスクノードでプロセスのサービス停止が成功したかどうかを判定する。成功したときは、処理をステップS07に進める。一部のディスクノードが、プロセスの停止に失敗、またはメタデータ変更が発生して、失敗(NG)応答を返してきたときは、処理をステップS09に進める。

【0076】

【ステップS07】 プロセスへのサービス停止が正常に終了したときは、続いて、新プロセスのサービス開始を指示する。指示を受けたディスクノードは、新プロセスによるサービスを開始させる。

40

【0077】

【ステップS08】 プロセスへの処理終了を指示する。指示を受けたディスクノードは、プロセスを終了させる。ディスクノードからの応答を受け、処理を終了する。

【ステップS09】 アップデート処理が途中で失敗、またはメタデータ変更が発生したときは、ロールバック処理を行って、アップデート開始前の状態に戻し、処理を終了する。

【0078】

ロールバック処理について説明する。

50

図13は、ロールバック処理の手順を示したフローチャートである。

【ステップS11】 失敗(NG)が、新プロセス起動時に発生したのかどうかを判定する。一部のディスクノードが新プロセス起動に失敗、または、メタデータ変更が発生したときは、処理をステップS15に進める。それ以外であれば、処理をステップS12に進める。

【0079】

【ステップS12】 失敗(NG)が、同期モードへの変更時に発生したのかどうかを判定する。一部のディスクノードがプロセスの同期化に失敗、または、メタデータ変更が発生したときは、処理をステップS14に進める。それ以外、ここでは、一部のディスクノードが、プロセスの停止に失敗、またはメタデータ変更が発生したときは、処理をステップS13に進める。

10

【0080】

【ステップS13】 サービス停止をしたプロセスのサービスを再開させる。

【ステップS14】 同期モードを非同期モードに戻す。

【ステップS15】 起動した新プロセスを終了させる。

【0081】

以上の処理手順が実行されることにより、アップデート処理が途中で失敗、またはアップデート処理の途中でメタデータが変更されたときは、実行された処理を逆順に辿り、アップデート開始前の状態に戻す。このように、アップデート開始前の状態に自動的に戻ることにより、サービス停止を抑止することができる。また、メタデータの変更が終了した時点などで、再び、アップデート処理を開始させることができる。

20

【0082】

なお、上記の処理機能は、コンピュータによって実現することができる。その場合、ストレージシステムを構成する管理ノード及びストレージノードが有すべき機能の処理内容を記述したプログラムが提供される。そのプログラムをコンピュータで実行することにより、上記処理機能がコンピュータ上で実現される。処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。

【0083】

プログラムを流通させる場合には、例えば、そのプログラムが記録されたDVD(Digital Versatile Disc)、CD-ROM(Compact Disc Read Only Memory)などの可搬型記録媒体が販売される。また、プログラムをサーバコンピュータの記憶装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピュータにそのプログラムを転送することもできる。

30

【0084】

プログラムを実行するコンピュータは、例えば、可搬型記録媒体に記録されたプログラムもしくはサーバコンピュータから転送されたプログラムを、自己の記憶装置に格納する。そして、コンピュータは、自己の記憶装置からプログラムを読み取り、プログラムに従った処理を実行する。なお、コンピュータは、可搬型記録媒体から直接プログラムを読み取り、そのプログラムに従った処理を実行することもできる。また、コンピュータは、サーバコンピュータからプログラムが転送されるごとに、逐次、受け取ったプログラムに従った処理を実行することもできる。

40

【0085】

以上の実施の形態に関し、更に以下の付記を開示する。

(付記1) 複数のストレージ装置に分散してデータを格納するストレージシステムにてコンピュータにアップデート処理を行わせるアップデート処理プログラムにおいて、

前記コンピュータを、

前記ストレージ装置へのデータ書き込み要求を受け取ったときに、書き込みデータを前記ストレージ装置に直接書き込む同期モードと、前記書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングで前記ストレージ装置に書き込む非同期モードとを備えた読み書き制御手段に対し、前記同期モードまたは前記非同期モードのどちらで

50

データ書き込みを行うかを指示する同期 / 非同期指示手段、

所定のサービス処理を実行しているプロセスと、新プロセスの動作を制御し、アップデートが要求されると、前記同期 / 非同期指示手段に指示して前記プロセスが前記サービス処理を実行中は前記非同期モードに設定される前記読み書き制御手段を前記同期モードに切り替え、前記プロセスの出力する前記書き込みデータが前記ストレージ装置に直接書き込まれている状態で前記プロセスに終了指示を出して前記サービス処理を終了させ、前記プロセスが前記サービス処理終了後に前記新プロセスのサービス処理を開始させるアップデート処理を行うプロセス制御手段、

として機能させることを特徴とするストレージシステムのアップデート処理プログラム

。

10

【0086】

(付記2) 前記プロセス制御手段は、前記新プロセスによるサービスを開始させるときは、前記同期 / 非同期指示手段に指示して、前記読み書き制御手段を前記非同期モードに切り替える、

ことを特徴とする付記1記載のストレージシステムのアップデート処理プログラム。

【0087】

(付記3) 前記ストレージ装置の記憶領域は、仮想的な論理ボリュームを所定の記憶領域単位で分割した論理セグメントと、前記論理セグメントに対応付けられた前記ストレージ装置の実データ記憶領域を前記所定の記憶領域単位で分割したスライスとを関連付けた管理情報によって管理されており、

20

前記プロセス制御手段は、前記同期モードを指示する前に、前記複数のストレージ装置に分散して格納される前記管理情報を読み出し、前記新プロセスを、前記サービス処理を停止させた状態で起動しておく、

ことを特徴とする付記1記載のストレージシステムのアップデート処理プログラム。

【0088】

(付記4) 前記コンピュータを、少なくとも、前記アップデート処理の開始前と、前記アップデート処理の終了時とに、前記ストレージ装置に分散して格納される、仮想的な論理ボリュームを所定の記憶領域単位で分割した論理セグメントと、前記論理セグメントに対応付けられた前記ストレージ装置の実データ記憶領域を前記所定の記憶領域単位で分割したスライスとを関連付けた管理情報を収集し、前記アップデート処理の開始前の前記管理情報と、前記アップデート処理の終了時の前記管理情報とを比較し、前記管理情報に変更があったか否かを前記プロセス制御手段に通知する管理情報チェック手段として機能させ、

30

前記プロセス制御手段は、前記管理情報の変更があったときは前記アップデート処理を中止し、前記アップデート処理が開始される前の状態に戻す、

ことを特徴とする付記1記載のストレージシステムのアップデート処理プログラム。

【0089】

(付記5) 前記プロセス制御手段は、前記複数のストレージ装置それぞれに接続される複数の前記コンピュータを管理する管理装置から前記コンピュータに向けて一斉に送信される指示に応じて、前記読み書き制御手段を前記同期モードで動作させる処理、前記プロセスに終了指示を出して処理を終了させる処理、及び前記新プロセスによるサービスを開始させる処理をそれぞれ実行し、処理が終了するごとに前記管理装置に終了結果を通知し、次の指示を待つ、

40

ことを特徴とする付記1記載のストレージシステムのアップデート処理プログラム。

【0090】

(付記6) 前記プロセス制御手段は、1または複数の前記コンピュータが前記指示に応じた処理に失敗したときに前記管理装置によって発行されるロールバック指示を受けたときは、前記管理装置からの指示に応じて実行してきた前記アップデート処理の手順を逆順に辿り、前記アップデート処理が開始される前の状態に戻す、

ことを特徴とする付記5記載のストレージシステムのアップデート処理プログラム。

50



## 【 0 0 9 1 】

(付記7) 前記コンピュータを、少なくとも、前記アップデート処理の開始前と、前記アップデート処理の終了時とに、前記ストレージ装置に分散して格納される、仮想的な論理ボリュームを所定の記憶領域単位で分割した論理セグメントと、前記論理セグメントに対応付けられた前記ストレージ装置の実データ記憶領域を前記所定の記憶領域単位で分割したスライスとを関連付けた管理情報を収集し、前記アップデート処理の開始前の前記管理情報と、前記アップデート処理の終了時の前記管理情報とを比較し、前記管理情報に変更があったか否かを前記プロセス制御手段に通知する管理情報チェック手段として機能させ、

前記プロセス制御手段は、前記管理情報チェック手段によって前記管理情報の変更が検出されたときは、前記管理装置に対し前記管理情報に変更があったことを通知する、ことを特徴とする付記6記載のストレージシステムのアップデート処理プログラム。

10

## 【 0 0 9 2 】

(付記8) 複数のストレージ装置に分散してデータを格納するストレージシステムのアップデート処理方法において、

所定のサービス処理を実行しているプロセスと、新プロセスの動作を制御するプロセス制御手段が、アップデートが要求されると、前記プロセスが前記サービス処理を実行中は前記ストレージ装置への書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングで前記ストレージ装置に書き込む非同期モードが選択される読み書き制御手段の動作モードを、前記書き込みデータを前記ストレージ装置に直接書き込む同期モードに切り替える指示を出す手順と、

20

同期/非同期指示手段が、前記プロセス制御手段の指示に従って、前記同期モードまたは前記非同期モードで前記プロセスまたは前記新プロセスの出力する前記書き込みデータを前記ストレージ装置に書き込む前記読み書き制御手段に対し、前記非同期モードから前記同期モードへ切り替える指示を出す手順と、

前記プロセス制御手段が、前記読み書き制御手段が前記同期モードで動作し、前記プロセスの前記書き込みデータが前記ストレージ装置に直接書き込まれている状態で前記プロセスに終了指示を出して前記サービス処理を終了させる手順と、

前記プロセス制御手段が、前記プロセスが前記サービス処理終了後に、前記新プロセスのサービス処理を開始させる手順と、

30

を有することを特徴とするストレージシステムのアップデート処理方法。

## 【 0 0 9 3 】

(付記9) 複数のストレージ装置に分散してデータを格納するストレージシステムにおいて、

前記ストレージ装置へのデータ書き込み要求を受け取ったときに、書き込みデータを前記ストレージ装置に直接書き込む同期モードと、前記書き込みデータをキャッシュメモリに蓄積し、蓄積データを所定のタイミングで前記ストレージ装置に書き込む非同期モードとを備えた読み書き制御手段に対し、前記同期モードまたは前記非同期モードのどちらでデータ書き込みを行うかを指示する同期/非同期指示手段と、所定のサービス処理を実行しているプロセスと、新プロセスの動作を制御し、入力された指示に応じて、前記同期/非同期指示手段に指示して前記プロセスが前記サービス処理を実行中は前記非同期モードに設定される前記読み書き制御手段を前記同期モードに切り替える処理、前記プロセスの出力する前記書き込みデータが前記ストレージ装置に直接書き込まれている状態で前記プロセスに終了指示を出して前記サービス処理を終了させる処理、及び前記プロセスが前記サービス処理終了後に前記新プロセスのサービス処理を開始させる処理を行い、処理が終了するごとに終了結果を通知するプロセス制御手段と、を備えたストレージノードと、

40

ネットワークを介して複数の前記ストレージノードに接続し、アップデート対象の前記プロセスを有する前記ストレージノードのアップデートの進行状況に関するアップデート管理情報を格納する管理情報記憶手段と、前記管理情報記憶手段に格納される前記アップデート管理情報に基づき、対象ストレージノードすべてに対し前記同期モードへの切り替

50

えを指示し、前記対象ストレージノードすべての処理が成功したときは前記対象ストレージノードに対し前記プロセスの停止を指示し、前記対象ストレージノードすべての処理が成功したときは前記対象ストレージノードに対し前記新プロセスを起動させ前記新プロセスによるサービスを開始させるアップデート管理手段と、を備えた管理ノードと、を有することを特徴とするストレージシステム。

【図面の簡単な説明】

【0094】

【図1】発明の概要を示した図である。

【図2】本実施の形態のマルチノードストレージの構成例を示す図である。

【図3】ディスクノードのハードウェア構成例を示す図である。

10

【図4】マルチノードストレージにおいてアップデート処理を行う各部のソフトウェア構成を示した図である。

【図5】管理テーブルの一例を示した図である。

【図6】コマンドの発行処理に応じた管理テーブルの変化を示した図である。

【図7】アップデート処理の動作シーケンス（同期モードへの切り替えまでの手順）を示した図である。

【図8】アップデート処理の動作シーケンス（新プロセスへの切り替えまでの手順）を示した図である。

【図9】新プロセス起動に失敗したときの動作シーケンスを示した図である。

【図10】プロセスの同期化に失敗したときの動作シーケンスを示した図である。

20

【図11】プロセスのサービス停止に失敗したときの動作シーケンスを示した図である。

【図12】管理ノードのアップデート処理手順を示したフローチャートである。

【図13】ロールバック処理の手順を示したフローチャートである。

【図14】従来のマルチノードストレージのプロセスアップデートの手順を示した図である。

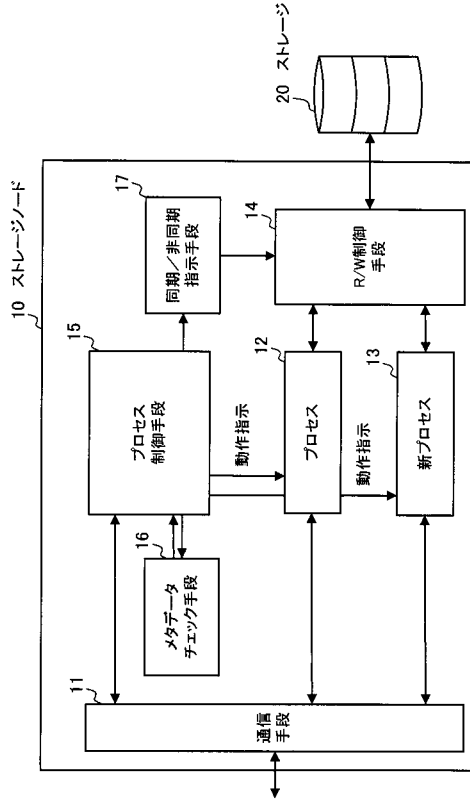
【符号の説明】

【0095】

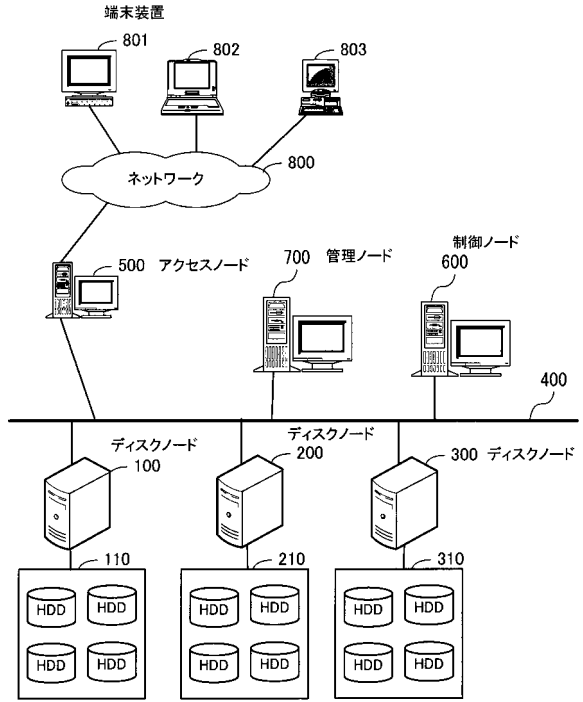
- 10 ストレージノード
- 11 通信手段
- 12 プロセス
- 13 新プロセス
- 14 読み書き（R/W）制御手段
- 15 プロセス制御手段
- 16 メタデータチェック手段
- 17 同期/非同期指示手段
- 20 ストレージ

30

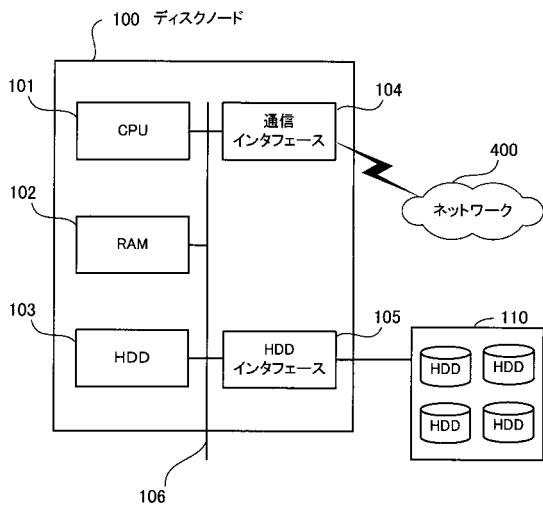
【図1】



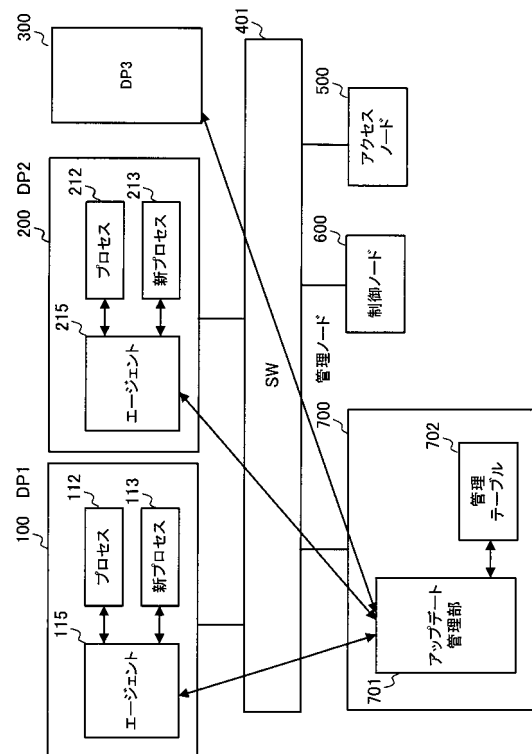
【図2】



【図3】



【図4】



【図5】

702 管理テーブル

ノードID	コマンド発行	結果
DP1	NULL/DONE	NULL/OK/NG
DP2	NULL/DONE	NULL/OK/NG
DP3	NULL/DONE	NULL/OK/NG

7021                      7022                      7023

【図6】

(A) コマンド発行前

ノードID	コマンド発行	結果
DP1	NULL	NULL
DP2	NULL	NULL
DP3	NULL	NULL

(B) コマンド発行後

ノードID	コマンド発行	結果
DP1	DONE	NULL
DP2	DONE	NULL
DP3	DONE	NULL

(C) コマンド発行後結果受け付け中

ノードID	コマンド発行	結果
DP1	DONE	OK
DP2	DONE	NULL
DP3	DONE	OK

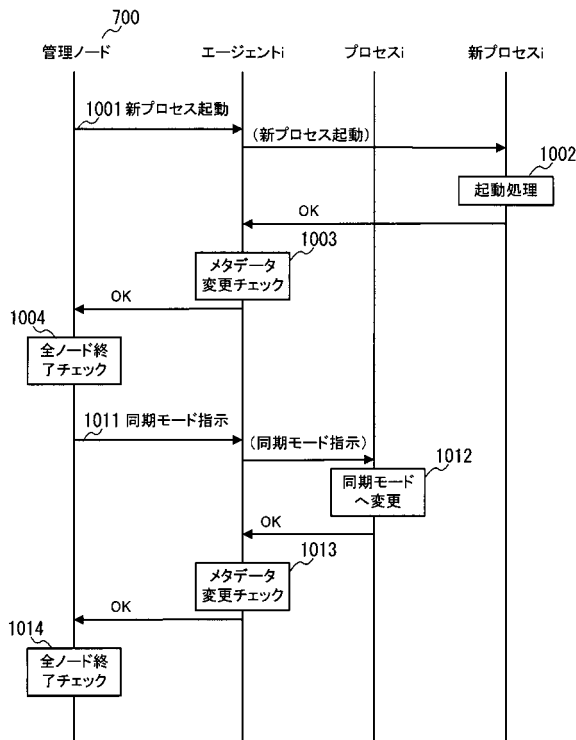
(D) コマンド結果受信後

ノードID	コマンド発行	結果
DP1	DONE	OK
DP2	DONE	OK
DP3	DONE	OK

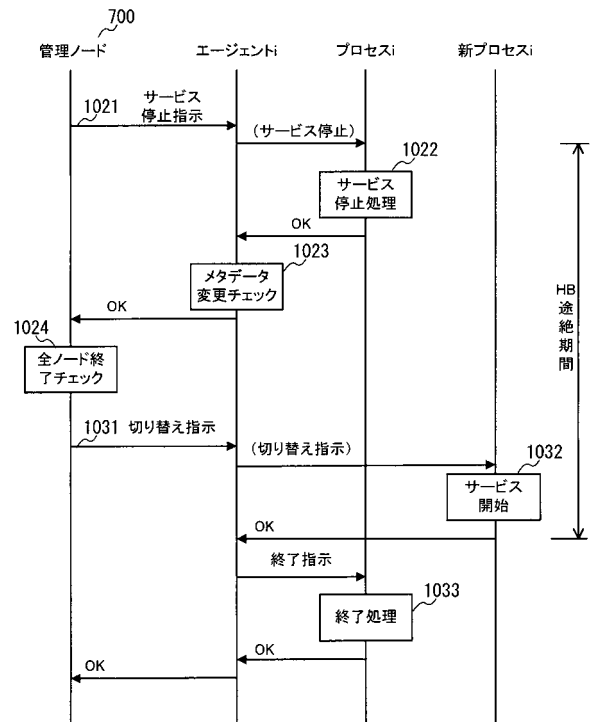
(E) コマンド結果受信後(NG含む)

ノードID	コマンド発行	結果
DP1	DONE	OK
DP2	DONE	NG
DP3	DONE	OK

【図7】

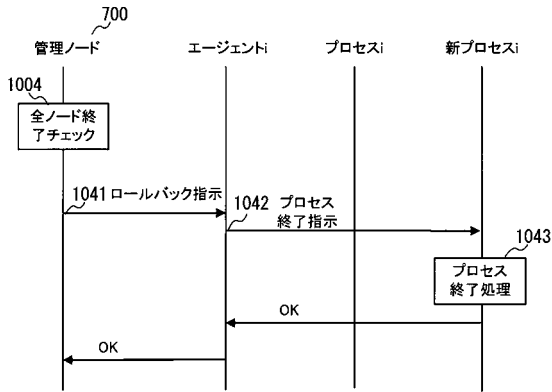


【図8】



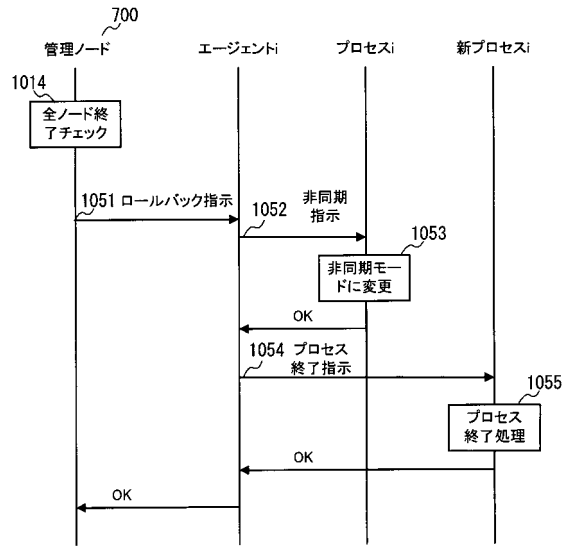
【図 9】

(a) 一部のノードで新プロセスの起動に失敗またはメタデータ変更発生



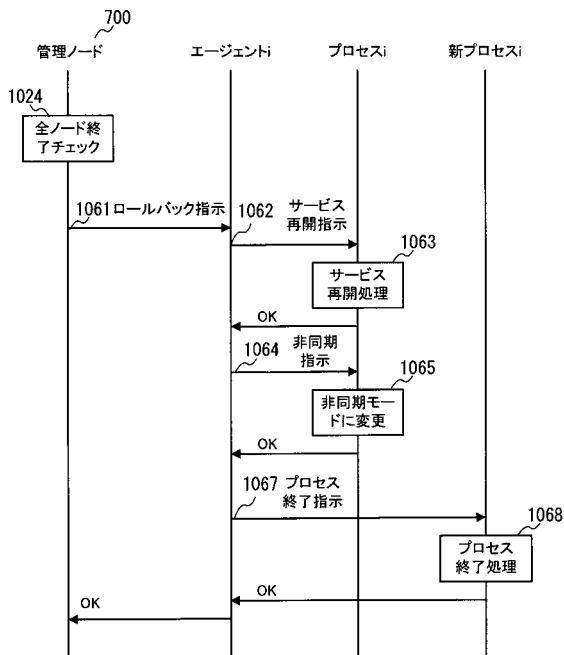
【図 10】

(b) 一部のノードでプロセスの同期化に失敗またはメタデータ変更発生

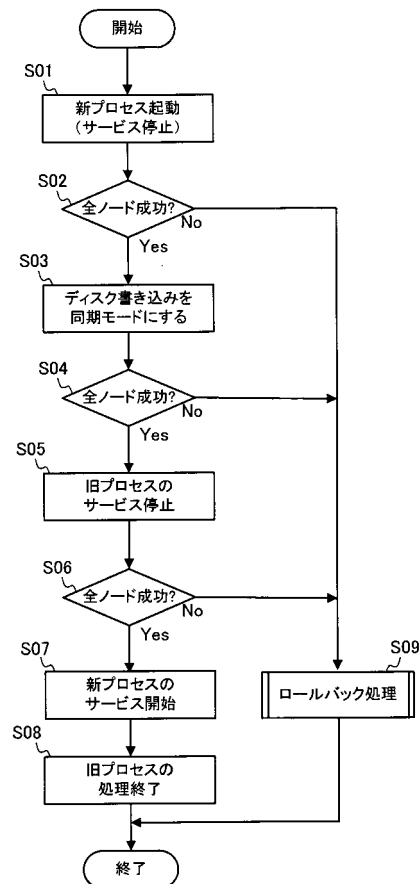


【図 11】

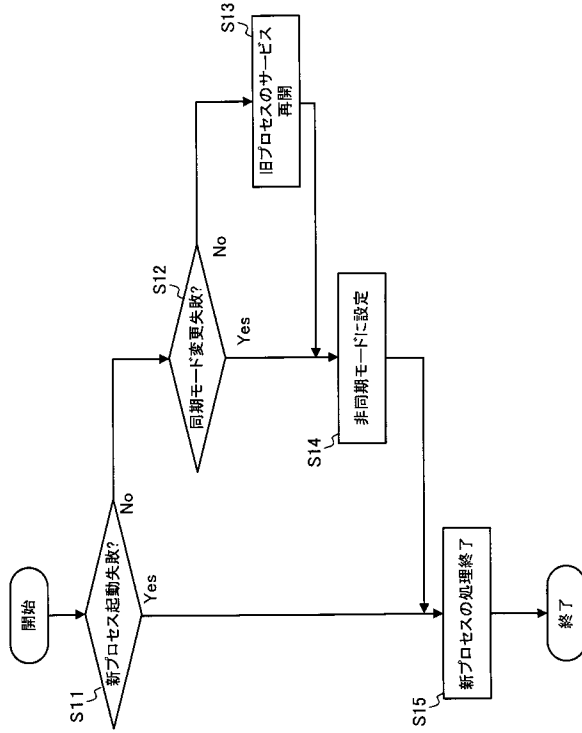
(c) 一部のノードでプロセスのサービス停止に失敗、またはメタデータ変更発生



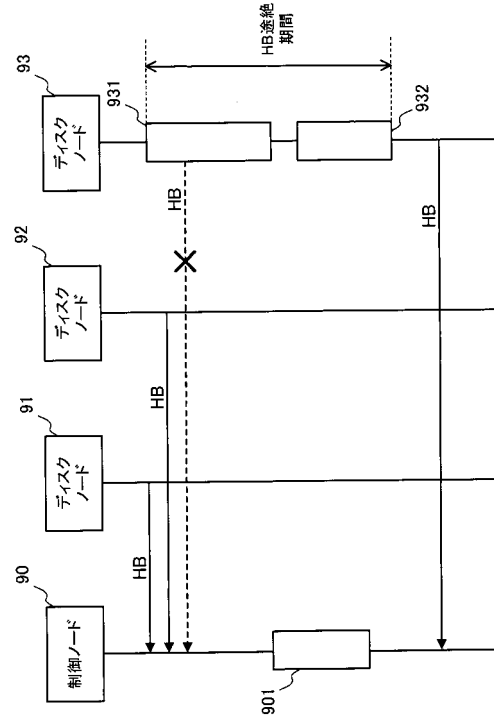
【図 12】



【図13】



【図14】



---

フロントページの続き

- (56)参考文献 特開2005 - 122763 (JP, A)  
特開2006 - 277205 (JP, A)  
特開2008 - 225763 (JP, A)  
特開2007 - 272496 (JP, A)  
特開2004 - 362589 (JP, A)  
特表2002 - 519765 (JP, A)  
国際公開第2008 / 136075 (WO, A1)  
特開平04 - 245352 (JP, A)  
米国特許出願公開第2005 / 0050267 (US, A1)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06