

【公報種別】特許法第17条の2の規定による補正の掲載

【部門区分】第6部門第3区分

【発行日】平成23年5月26日(2011.5.26)

【公表番号】特表2003-515813(P2003-515813A)

【公表日】平成15年5月7日(2003.5.7)

【出願番号】特願2001-540586(P2001-540586)

【国際特許分類】

G 06 F	15/177	(2006.01)
G 06 F	9/46	(2006.01)
G 06 F	12/00	(2006.01)

【F I】

G 06 F	15/177	6 8 2 G
G 06 F	9/46	3 6 0 E
G 06 F	12/00	5 0 1 A
G 06 F	12/00	5 1 4 E

【誤訳訂正書】

【提出日】平成23年4月7日(2011.4.7)

【誤訳訂正1】

【訂正対象書類名】明細書

【訂正対象項目名】全文

【訂正方法】変更

【訂正の内容】

【書類名】明細書

【発明の名称】ストレージネットワーク内のクオーラムリソースアービタ

【特許請求の範囲】

【請求項1】複数のコンピュータノードと、1つまたは複数の記憶デバイス内に含まれる物理クオーラムリソースを具備する記憶サブシステムとを備えるシステムにおいて、前記物理クオーラムリソースがストレージネットワークの前記コンピュータノードに利用可能である前記ストレージネットワークを形成する方法であって、

前記システムの前記コンピュータノードのそれぞれにより前記記憶サブシステムの前記物理クオーラムリソースについての何れかの現在の所有権をプロセッサが終了させるステップと、

新しいストレージネットワークの構成情報を更新するために前記システムのそれぞれの他のコンピュータノードを前記プロセッサがスキャンするステップと、

1つまたは複数の前記コンピュータノードにより前記物理クオーラムリソースの所有権を決定するためにアービトレーションを前記プロセッサが呼び出すステップと、

前記物理クオーラムリソースを備えるクオーラムボリュームを前記ストレージネットワークのメンバーに前記プロセッサが組み込むステップとを備え、

前記現在の所有権を終了させるステップと、前記アービトレーションを呼び出すステップは、第1のプログラムモジュールにより実行され、前記それぞれの他のコンピュータノードをスキャンするステップと、前記クオーラムボリュームを組み込むステップは、前記第1のプログラムモジュールと異なる第2のプログラムモジュールにより実行されることを特徴とする方法。

【請求項2】前記現在の所有権を終了させるステップは、前記記憶サブシステムへのアクセスをロックすることを含むことを特徴とする請求項1に記載の方法。

【請求項3】前記現在の所有権を終了させるステップは、それぞれのコンピュータノードのバスをリセットすることを含むことを特徴とする請求項1に記載の方法。

【請求項4】前記現在の所有権を終了させるステップは、アービトレーションを呼

び出す前に所定の遅延期間待機することを含むことを特徴とする請求項 1 に記載の方法。

【請求項 5】 前記ストレージネットワークの前記コンピュータノードにより所有された物理クオーラムリソースを含む前記記憶サブシステムのそれぞれの記憶デバイスを識別する構成情報を前記プロセッサが生成するステップを更に備えることを特徴とする請求項 1 に記載の方法。

【請求項 6】 内部構成データベースを再構築するために、前記生成された構成情報に含まれるそれぞれの記憶デバイスからの前記クオーラムボリュームの情報を前記プロセッサが処理するステップを更に備えることを特徴とする請求項 5 に記載の方法。

【請求項 7】 前記クオーラムボリュームを組み込むステップは、コンピュータノードがクオーラムボリュームに必要な全ての記憶デバイスの所有権を獲得したとき実行されることを特徴とする請求項 1 に記載の方法。

【請求項 8】 前記クオーラムボリュームを組み込むステップは、コンピュータノードが前記物理クオーラムリソースの過半数の所有権を獲得したとき実行されることを特徴とする請求項 1 に記載の方法。

【請求項 9】 1 以上の前記物理クオーラムリソースは、1 つまたは複数の連結されかつストライプされたボリュームエクステントを具備し、

前記クオーラムボリュームを組み込むステップは、所有権が 1 より大きい前記連結されかつストライプされたボリュームエクステントのために獲得されたとき実行されることを特徴とする請求項 1 に記載の方法。

【請求項 10】 1 以上の前記物理クオーラムリソースは、1 つまたは複数の連結されかつストライプされたボリュームエクステントを具備し、

前記クオーラムボリュームを組み込むステップは、所有権が前記 1 つまたは複数の連結されかつストライプされたボリュームエクステントの単純な過半数のために獲得されたとき実行されることを特徴とする請求項 1 に記載の方法。

【請求項 11】 複数のコンピュータノードと、1 つまたは複数の記憶デバイス内に含まれる物理クオーラムリソースを具備する記憶サブシステムとを備えるシステム内で、前記物理クオーラムリソースがストレージネットワークの前記コンピュータノードに利用可能である前記ストレージネットワークを形成するための各ステップを実施するコンピュータ実行可能命令を記憶するコンピュータ可読記録媒体であって、前記各ステップは、

前記システムの前記コンピュータノードのそれぞれにより前記記憶サブシステムの前記物理クオーラムリソースについての何れかの現在の所有権をプロセッサが終了させるステップと、

新しいストレージネットワークの構成情報を更新するために前記システムのそれぞれの他のコンピュータノードを前記プロセッサがスキャンするステップと、

1 つまたは複数の前記コンピュータノードにより前記物理クオーラムリソースの所有権を決定するためにアービトレーションを前記プロセッサが呼び出すステップと、

前記物理クオーラムリソースを備えるクオーラムボリュームを前記ストレージネットワークのメンバーに前記プロセッサが組み込むステップとを備え、

前記現在の所有権を終了させるステップと、前記アービトレーションを呼び出すステップは、第 1 のプログラムモジュールにより実行され、前記それぞれの他のコンピュータノードをスキャンするステップと、前記クオーラムボリュームを組み込むステップは、前記第 1 のプログラムモジュールと異なる第 2 のプログラムモジュールにより実行されることを特徴とするコンピュータ可読記録媒体。

【請求項 12】 前記現在の所有権を終了させるステップは、前記記憶サブシステムへのアクセスをロックすることを含むことを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 13】 前記現在の所有権を終了させるステップは、それぞれのコンピュータノードのバスをリセットすることを含むことを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 14】 前記現在の所有権を終了させるステップは、アービトレーションを

呼び出す前に所定の遅延期間待機することを含むことを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 15】 前記ストレージネットワークの前記コンピュータノードにより所有された物理クオーラムリソースを含む前記記憶サブシステムのそれぞれの記憶デバイスを識別する構成情報を前記プロセッサが生成するステップを更に備えることを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 16】 内部構成データベースを再構築するために、前記生成された構成情報に含まれるそれぞれの記憶デバイスからの前記クオーラムボリュームの情報を前記プロセッサが処理するステップを更に備えることを特徴とする請求項 15 に記載のコンピュータ可読記録媒体。

【請求項 17】 前記クオーラムボリュームを組み込むステップは、コンピュータノードがクオーラムボリュームに必要な全ての記憶デバイスの所有権を獲得したとき実行されることを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 18】 前記クオーラムボリュームを組み込むステップは、コンピュータノードが前記物理クオーラムリソースの過半数の所有権を獲得したとき実行されることを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 19】 1 以上の前記物理クオーラムリソースは、1つまたは複数の連結されかつストライプされたボリュームエクステントを具備し、

前記クオーラムボリュームを組み込むステップは、所有権が1より大きい前記連結されかつストライプされたボリュームエクステントのために獲得されたとき実行されることを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【請求項 20】 1 以上の前記物理クオーラムリソースは、1つまたは複数の連結されかつストライプされたボリュームエクステントを具備し、

前記クオーラムボリュームを組み込むステップは、所有権が前記1つまたは複数の連結されかつストライプされたボリュームエクステントの単純な過半数のために獲得されたとき実行されることを特徴とする請求項 11 に記載のコンピュータ可読記録媒体。

【発明の詳細な説明】

【0001】

(関連出願)

本出願は、本出願と同日に出願され同じ譲受人に譲渡された、以下のすべての出願に関連するものである。

第 09 / 449577 号(整理番号 777245US1) 「Storage Management System Having Common Volume Manager」

第 09 / 450364 号(整理番号 777246US1) 「Storage Management System Having Abstracted Volume Providers」

第 09 / 451219 号(整理番号 777247US1) 「Volume Stacking Model」

第 09 / 450300 号(整理番号 777248US1) 「Volume Configuration Data Administration」

第 0 / 451220 号(整理番号 777249US1) 「Volume Migration Between Volume Groups」

【0002】

(発明の分野)

本発明は、一般にデータ記憶デバイスに関し、より詳細には、ストレージネットワーク内での論理クオーラムリソースのアビトリエーション機構に関する。

【0003】

(版権の通知 / 許可)

本特許文書の開示の一部には、版権保護を対象とする資料が含まれている。版権所有者

は、特許商標事務所の特許ファイルまたは記録に表示されているものと同じ特許文書または特許開示を複製することに異議はないが、そうでない場合はあらゆる版権を保有する。以下の通知は、下記および図面に記載するようなソフトウェアおよびデータに適用される。

Copyright 1999, Microsoft Corporation,
All Rights Reserved.

【0004】

(発明の背景)

コンピュータシステムを発展させる場合、磁気ディスクや光ディスクなどのデータ記憶デバイスを利用および構成することができる。たとえば、これらの記憶デバイスは、バスを介してコンピュータシステムに接続するか、または無線または有線のネットワークを介してコンピュータシステムに接続することができる。さらにこれらの記憶デバイスは、分離するかまたは单一キャビネット内に配置することができる。

【0005】

ストレージネットワークとは、ノードと呼ばれ、単一の記憶資源として動作する、相互接続されたコンピューティングシステムの集まりのことである。ストレージネットワークでは、ハードウェアまたはソフトウェアに障害が発生した場合でもシステムの動作を続行させることができあり、ノードを簡単に追加できるようにすることでスケーラビリティを高くし、管理者がノードを単一のシステムとして管理できるようにすることで管理を簡略化する。

【0006】

クラスタソフトウェアはそれぞれのノード上にあり、ストレージネットワークのクラスタに特有のすべての活動を管理する。クラスタソフトウェアは、しばしば、ノードのスタートアップを自動的に実行する。この際、クラスタソフトウェアはローカルの非共有デバイスを構成して装着する。また、クラスタソフトウェアは、ストレージネットワークの他のメンバが動作可能であるかどうかを判定するための「発見」プロセスも使用する。クラスタソフトウェアは既存のクラスタを発見すると、認証シーケンスを実行することによってそのクラスタの接合を試みる。既存クラスタのクラスタマスターは新規クラスタを認証し、接合ノードの認証が成功すれば、その状況を戻す。ノードがメンバとして認識されない場合は、接合要求が拒否される。

【0007】

発見プロセス中にクラスタが見つからない場合、ノードは独自のクラスタの形成を試みる。このプロセスは、ノードが自分の属するクラスタと通信できないときには何回でも繰り返される。従来のコンピューティングシステムでは、ノードは、ストレージネットワークを形成するために、ディスクなどの物理「クオーラムリソース」(physical "quorum resource")に対してアービトレイションを実行する。最近のシステムでは、クオーラムリソースは、1つまたは複数の物理クオーラムリソースを含むボリュームなどの論理リソースが可能である。たとえばボリュームは、ディスクの一部分、ディスク全体、複数ディスクの一部分、または複数ディスクである可能性のある論理記憶ユニットである。

【0008】

従来のシステムでは、クラスタの所有権を決定する、すなわちアービトレイションプロセスを実施するための役割および情報が、いくつかの構成要素および/またはソフトウェアモジュール間で分配されることが多い。基礎となる記憶デバイスを構成および管理するための役割も、同様に分配されることが多い。このように役割を明確に分割しないと、所与の構成要素またはソフトウェアモジュールを変更することが困難になる。したがって、当分野では、クラスタ管理からのクラスタアービトレイションの役割と、ボリューム管理および基礎となる記憶デバイスの役割とを、より明確に分離することが求められている。

【0009】

(発明の概要)

本発明は、上記の短所、欠点、および問題点に対処するものである。本発明のクラスタ管理ソフトウェアおよびボリューム管理ソフトウェアは、ストレージネットワークのノード上で実行され、クオーラムボリュームなどの論理クオーラムリソースをアビトリエーションするために、基礎となるオペレーティングシステムと協働して動作する。本発明によれば、クラスタ管理ソフトウェアは論理クオーラムリソースをアビトリエーションし、基礎となる物理クオーラムリソースの知識を持たずにストレージネットワークを形成する。この方法では、クラスタ管理ソフトウェアはハードウェア特有のものではない。さらに、クラスタ管理ソフトウェアは、基礎となる物理クオーラムリソースから論理クオーラムリソースがどのように形成されるかを知っておく必要がない。たとえば、ボリューム管理ソフトウェアが管理する役割は、論理クオーラムボリュームの形成および装着のみである。ボリューム管理ソフトウェアは、アビトリエーションプロセスおよび所有権決定に関する詳細な知識を持たずに、ボリューム管理を実行する。

【0010】

(発明の詳細な説明)

本発明の例示的な実施形態を詳細に説明する以下の記述では、本明細書の一部であり、本発明が実施可能な特有の例示的な実施形態が例として示されている、添付の図面を参照する。これらの実施形態は、当分野の技術者が本発明を実施できるように詳細に記述されており、他の実施形態が利用可能であること、および本発明の範囲を逸脱することなく変更が可能であることが理解されよう。したがって、以下の詳細な説明は限定的なものではなく、本発明の範囲は特許請求の範囲によってのみ定義される。

【0011】

詳細な説明は、4つのセクションに分けられる。第1のセクションは用語集である。第2のセクションでは、本発明の実施形態がそれと共に実施可能なハードウェアおよび動作環境について説明する。第3のセクションでは、本発明のシステムレベルの概要を示す。最後に第4のセクションでは、詳細な説明の結果を示す。

【0012】

[用語の定義]

損失：フォールトトレラントボリュームに、1つまたは複数のディスクまたはボリュームエクステントがないことを示す状況。たとえば、現在1つのミラーしか使用できないミラーセット。

【0013】

構成データ：物理リソースを論理ボリュームにマッピングすることを示す。

【0014】

指示された構成：プロバイダに、論理ブロックの再マッピングを選択するための規則が明示的に与えられること。

【0015】

ディスクプラッタ：ディスクパックからボリュームをエクスポートまたはインポートするのに使用される、ディスクパックのサブセット。

【0016】

ディスクパック：論理ボリュームと基礎となるディスクの集まり。ディスクパックとは、ボリュームの推移閉包の単位である。

【0017】

エクスポート：ディスクプラッタおよびそのプラッタに含まれるすべてのボリュームを、1つのディスクパックから移動すること。

【0018】

露出された：ボリュームに、関連するボリューム名（ドライブ文字）または装着位置がある場合、ボリュームはオペレーティングシステムに対して露出されている。ボリュームは、ファイルシステムまたは他のデータ格納に使用可能にすることができる。

【0019】

フリー エージェント ドライブ：ディスクパックのメンバでないディスク ドライブ。フリー エージェント ドライブは、露出された論理ボリュームを含むことはできない。

【0020】

健全：ボリューム障害管理状況。ボリュームの状況には、初期、健全、損失、不健全、または再構築がある。

【0021】

健全な：有効データを含んでいるかまたは含むことができる。

【0022】

ホットスポットティング：ボリュームまたはボリュームエクステントの集まりを一時的にプレクシング（p l e x i n g）すること。

【0023】

インポート：ディスクプラッタおよびそのプラッタに含まれるすべてのボリュームを、1つのディスクパックに移動すること。

【0024】

初期化：ボリュームがボリューム構成を再発見していることを示す状況。

【0025】

LBN：論理ブロック番号

【0026】

論理ブロックマッピング：論理ボリュームプロバイダに対して露出された論理ブロックと、同じプロバイダによって露出された論理ブロックとの関係。

【0027】

論理クオーラムリソース：ストレージネットワークの形成に必要な論理リソース。論理ボリュームなどの論理クオーラムリソースには、ディスクなどの1つまたは複数の物理クオーラムリソースが含まれる。

【0028】

論理ボリューム：ディスクの一部分、ディスク全体、複数ディスクの一部分、または複数ディスクである可能性のある、論理記憶ユニット。

【0029】

論理ボリュームプロバイダ：論理ボリュームを露出するソフトウェア。プロバイダには、ランタイムサービス、構成データ、および管理サービスが含まれる。

【0030】

管理サービス：ボリュームの構成、モニタリング、または障害処理を実行するために、まれにしか実行されないソフトウェア。

【0031】

マッピング済みボリューム：単一のより大きなボリュームを露出するためにボリュームを連結する、単純な線形の論理ブロックマッピング。

【0032】

ミラーリング済みボリューム：2つまたはそれ以上の同一のデータコピーを維持する論理ボリューム。RAID 1とも呼ばれる。

【0033】

パーティストライピング済みボリューム：パーティティチェック情報ならびにデータを維持する論理ボリューム。厳密なマッピングおよび保護方式は、ベンダ特有のものである。RAID 3、4、5、6が含まれる。

【0034】

プレクシング済みボリューム：動的なミラー ボリューム。プレクシングは、フォールトトレランスを提供するのではなく、ボリュームのコピーを作成するのに使用される。ミラーは、コンテンツが同期化された後に除去する意図で、ボリュームに追加される。

【0035】

RAID：Redundant Array of Independent Disks

【 0 0 3 6 】

再構築：以前に損失したフォールトトレラントボリュームがすべてのボリュームエクステントデータを再同期化することを示す状況。

【 0 0 3 7 】

ランタイムサービス：I/Oごとの要求ベースで実行されるソフトウェア。

【 0 0 3 8 】

S C S I : S m a l l - C o m p u t e r S y s t e m s I n t e r f a c e

【 0 0 3 9 】

スタック済みボリューム：複数の論理ブロックマッピング動作によって構築されたボリューム。一例として、ミラーボリュームのストライプセットが挙げられる。スタッキングには、ストリッピング、マッピング、およびプレクシングが含まれる。

【 0 0 4 0 】

ストライピング済みボリューム：連続する論理ボリュームエクステントを複数のボリュームにまたがって分配する、論理ブロックマッピング。R A I D 0とも呼ばれる。

【 0 0 4 1 】

不健全：非フォールトトレラントボリュームに、1つまたは複数のディスクまたはボリュームエクステントがないことを示す状況。不健全ボリュームに含まれるデータにアクセスしてはいけない。

【 0 0 4 2 】

ボリューム構成安定性：ボリュームの論理対物理マッピングが変化しているかどうか。ボリュームには、安定、拡張、縮小、プレクシング、または再マッピングがある。

【 0 0 4 3 】

ボリュームエクステント：ボリューム上に含まれる論理ブロックの連続する領域。ボリュームエクステントは、最小管理対象論理ボリューム単位である。

【 0 0 4 4 】

ボリューム状況：システムによる現在のボリューム使用状況。ボリュームには、未使用、ホットスペア、マップ済み、使用、または未知がある。

【 0 0 4 5 】**[ハードウェアおよび動作環境]**

図1は、本発明の実施形態がそれと共に実施可能なハードウェアおよび動作環境を示す図である。図1の説明は、本発明がそれと共に実施可能な、適切なコンピュータハードウェアおよび適切なコンピューティング環境を、簡単かつ概略的に説明することを意図している。必須ではないが、本発明は、パーソナルコンピュータなどのコンピュータによって実行される、プログラムモジュールなどのコンピュータ実行可能命令の一般的な情況で説明される。一般に、プログラムモジュールには、特定のタスクを実行するか、または、特定の抽象データ型を実施する、ルーチン、プログラム、オブジェクト、構成要素、データ構造などが含まれる。

【 0 0 4 6 】

さらに当分野の技術者であれば、本発明が、ハンドヘルドデバイス、マルチプロセッサシステム、マイクロプロセッサベースまたはプログラム可能な大衆消費電子製品、ネットワークP C、ミニコンピュータ、メインフレームコンピュータなどを含む、他のコンピュータシステム構成でも実施可能であることを理解されよう。本発明は、通信ネットワークを介してリンクされたリモート処理デバイスによってタスクが実行される、分散型コンピューティング環境でも実施可能である。分散型コンピューティング環境では、プログラムモジュールは、ローカルおよびリモートのどちらのメモリ記憶デバイスにも配置可能である。

【 0 0 4 7 】

本発明を実施するための図1の例示的なハードウェアおよび動作環境には、コンピュータ20の形式の汎用コンピューティングデバイスが含まれ、これには、処理ユニット21、システムメモリ22、ならびに、システムメモリ22を含む様々なシステム構成要素を

処理ユニット21に動作可能に結合するシステムバス23が含まれる。処理ユニット21は1つだけでも複数であってもよく、その結果、コンピュータ20のプロセッサには、単一の中央処理ユニット(CPU)が含まれるか、または一般に並列処理環境と呼ばれる複数の処理ユニットが含まれる。コンピュータ20は、従来のコンピュータ、分散型コンピュータ、または任意の他のタイプのコンピュータが可能であり、本発明はそのように限定されていない。

【0048】

システムバス23は、任意の様々なバスアーキテクチャを使用する、メモリバスまたはメモリ制御装置、周辺バス、およびローカルバスを含む、いくつかのタイプのバス構造のうちいずれかでよい。システムメモリは単にメモリと呼ぶことも可能であり、読み取り専用メモリ(ROM)24およびランダムアクセスメモリ(RAM)25を含む。起動中などにコンピュータ20内の要素間で情報を転送するのを助ける基本ルーチンを含む基本入出力システム(BIOS)26が、ROM24に格納される。さらにコンピュータ20は、図示されていないハードディスクからの読み取りおよびハードディスクへの書き込みのためのハードディスクドライブ27、取外し可能磁気ディスク29からの読み取りまたはこのディスクへの書き込みのための磁気ディスクドライブ28、ならびに、CD-ROMまたは他の光学式媒体などの取外し可能光ディスク31からの読み取りまたはこのディスクへの書き込みのための光ディスクドライブ30を含む。

【0049】

ハードディスクドライブ27、磁気ディスクドライブ28、および光ディスクドライブ30は、それぞれハードディスクドライブインターフェース32、磁気ディスクドライブインターフェース33、および光ディスクドライブインターフェース34によって、システムバス23に接続される。ドライブおよびこれらに関連付けられたコンピュータ可読媒体は、コンピュータ可読命令、データ構造、プログラムモジュール、およびコンピュータ20に関する他のデータの、不揮発性記憶域を提供する。当分野の技術者であれば、磁気カセット、フラッシュメモリカード、デジタルビデオディスク、ベルヌイカートリッジ、ランダムアクセスメモリ(RAMs)、読み取り専用メモリ(ROMs)などの、コンピュータによってアクセス可能なデータを格納できる任意のタイプのコンピュータ可読媒体が、例示的動作環境で使用可能であることを理解されたい。

【0050】

いくつかのプログラムモジュールが、ハードディスク27、磁気ディスク29、光ディスク31、ROM24、あるいは、オペレーティングシステム35、1つまたは複数のアプリケーションプログラム36、他のプログラムモジュール37、およびプログラムデータ38を含むRAM25上に格納可能である。ユーザは、コマンドおよび情報を、キーボード40およびポインティングデバイス42などの入力デバイスを介して、パーソナルコンピュータ20に入力することができる。他の入力デバイス(図示せず)には、マイクロフォン、ジョイスティック、ゲームパッド、衛星放送用アンテナ、スキャナなどが含まれる。これらおよび他の入力デバイスは、システムバスに結合されたシリアルポートインターフェース46を介して、処理ユニット21に接続されることが多いが、パラレルポート、ゲームポート、またはユニバーサルシリアルバス(USB)などの他のインターフェースによって接続されることも可能である。モニタ47または他のタイプの表示デバイスも、ビデオアダプタ48などのインターフェースを介して、システムバス23に接続される。コンピュータには、モニタに加えて、典型的にはスピーカおよびプリンタなどの他の周辺出力デバイス(図示せず)が含まれる。

【0051】

コンピュータ20は、リモートコンピュータ49などの1つまたは複数のリモートコンピュータへの論理接続を使用して、ネットワーク化された環境で動作することができる。これらの論理接続は、ローカルコンピュータであるコンピュータ20に結合されたかまたはこの一部である通信デバイスによって達成され、本発明は特定タイプの通信デバイスに限定されるものではない。リモートコンピュータ49は、他のコンピュータ、サーバ、ル

ータ、ネットワークPC、クライアント、ピア(peer)デバイス、または他の共通ネットワークノードであってよく、典型的には、コンピュータ20に関連した前述の要素の多くまたはすべてを含むが、図1にはメモリ記憶デバイス50のみが図示されている。図1に示されている論理接続には、ローカルエリアネットワーク(LAN)51およびワイドエリアネットワーク(WAN)52が含まれる。こうしたネットワーキング環境は、事務所、全社的なコンピュータネットワーク、インターネット、およびインターネットで、一般的に見られるものである。

【0052】

コンピュータ20は、LANネットワーキング環境で使用される場合、通信デバイスの一種であるネットワークインターフェースまたはアダプタ53を介して、ローカルネットワーク51に接続される。コンピュータ20は、WANネットワーキング環境で使用される場合、典型的には、モデム54、通信デバイスの一種、あるいはインターネットなどのワイドエリアネットワーク52を介して通信を確立するための他の任意のタイプの通信デバイスを含む。モデム54は内部または外部であってよく、シリアルポートインターフェース46を介してシステムバス23に接続される。ネットワーク化された環境では、パソコン用コンピュータ20に関連して、またはその一部として示されたプログラムモジュールは、リモートメモリ記憶デバイスに格納することができる。図示されたネットワーク接続は例示的なものであり、コンピュータ間に通信リンクを確立するための他の手段および通信デバイスが使用可能であることが理解されよう。

【0053】

以上、本発明の実施形態がそれと共に実施可能なハードウェアおよび動作環境について説明してきた。本発明の実施形態がそれと共に実施可能なコンピュータは従来のコンピュータ、分散型コンピュータ、または任意の他のタイプのコンピュータであってよく、本発明はそのように限定されていない。こうしたコンピュータには、典型的にはその処理装置として1つまたは複数の処理ユニット、およびメモリなどのコンピュータ可読媒体が含まれる。コンピュータは、ネットワークアダプタまたはモデムなどの通信デバイスも含むことができるため、他のコンピュータと通信上で結合することができる。

【0054】

[システムレベルの概要]

図2は、ネットワーク120を介してノード110に通信上で結合されたノード105を含む、ストレージネットワーク100のシステムレベルの概要を示す構成図である。ノード105および110は、図1に示されたローカルコンピュータ20またはリモートコンピュータ49などの、任意の適切なコンピューティングシステムを表す。

【0055】

ストレージネットワーク100は、記憶デバイス107、記憶デバイス108、および記憶デバイス109を含む、記憶サブシステム106をさらに含む。これらのデバイスは、単一の内部ディスク、複数の外部ディスク、またはRAIDキャビネットなどの、任意の適切な記憶媒体であってよい。記憶サブシステム106は、デュアル接続SCSI(Small-Computer Systems Interface)、ファイバチャネルなどの任意の適切な相互接続メカニズムである、バス112を介して結合される。

【0056】

ストレージネットワーク100を形成するために、ノード105および110は、クオーラムボリュームなどの論理クオーラムリソースをアービトレーションする。図2では、論理クオーラムリソースが、物理クオーラムリソース111によって集合的に形成されたクオーラムボリュームとして示されており、本実施形態では、データ記憶デバイス108およびデータ記憶デバイス109内のデータ記憶エクステントである。ノード105または110のいずれかがすべての物理クオーラムリソース111の所有権を得ることに成功すると、成功したノードはストレージネットワーク100を形成することができる。以下で説明するように、本発明のクラスタ管理ソフトウェアおよびボリューム管理ソフトウェアは各ノード上で実行し、物理クオーラムリソース111の所有権がノード105と11

0とで分けられている状況を解決する。各ノード上で、クラスタ管理ソフトウェアおよびボリューム管理ソフトウェアが、ストレージネットワーク100を形成するために基礎となるオペレーティングシステムと協働する。以下で説明するように、アービトレーションおよび管理の役割は、クラスタ管理ソフトウェアとボリューム管理ソフトウェアで分担され、その結果、クラスタ管理ソフトウェアは、ボリューム管理の詳細および記憶サブシステム106について知らずに、アービトレーションプロセスを処理する。ボリューム管理ソフトウェアは、ストレージネットワーク100がどのように形成されるかを知らずに、記憶サブシステム106の構成および管理を処理する。

【0057】

図3は、様々な協働するソフトウェア構成要素が本発明のアービトレーション技法を実行する、図2のノード105またはノード110などのノード200の一実施形態を示す構成図である。ノード200では、クラスタマネージャ202がすべてのクラスタに特有の活動を監視し、ディスク制御装置206を介して記憶サブシステム106のバス112(図2)と通信する。クラスタマスターとして、クラスタマネージャ202、ボリュームマネージャ204、およびオペレーティングシステム35が協働して、ストレージネットワーク100および対応する物理クオーラムリソース111のクオーラムボリュームを管理する。より具体的に言えば、クラスタマネージャ202は、ボリューム管理および記憶サブシステム106の詳細を知らずに、アービトレーションプロセスを処理する。ボリュームマネージャは、すべてのボリュームマッピングおよびストレージネットワーク100の記憶サブシステム106の構成を処理する。ディスク制御装置206は、記憶サブシステム106とのすべての通信を処理し、SCSI、IPなどの様々なデータ通信プロトコルのうち1つを実施することができる。アプリケーション210は、ストレージネットワーク100と対話する、任意のユーザモードソフトウェアモジュールを表す。以上、本セクションでは、本発明の例示的な実施形態のシステムレベルの動作概要について、詳細に説明した。

【0058】

(本発明の例示的な実施形態の方法)

前のセクションでは、本発明の例示的な実施形態のシステムレベルの動作概要について説明した。このセクションでは、例示的な実施形態を実行中のコンピュータによって実行される特定の方法について、一連の流れ図を参照しながら説明する。コンピュータによって実行される方法は、コンピュータ実行可能命令で構成されるコンピュータプログラムとなる。流れ図を参照しながらこの方法について説明することで、当分野の技術者は、適切なコンピュータ(コンピュータ読み取り可能媒体からの命令を実行中のコンピュータのプロセッサ)上で方法を実行するためのこのような命令を含む、このようなプログラムを開くことができる。

【0059】

図4は、本発明が、クラスタ管理の役割とボリューム管理の役割を明確に分離する方法を示す図である。より具体的に言えば、アービトレーションサイクル300は、ストレージネットワーク100の各ノード上でクラスタマネージャ202およびボリュームマネージャによって実行されるときの、本発明の変形アービトレーション方法の一実施形態を示す。アービトレーションサイクル300は、ストレージネットワーク100がまだ確立されていないとき、たとえば、ノード105または110のいずれかが最初にブートされるときか、またはストレージネットワーク100があらかじめ形成されているがノード105と110の間の通信が故障しているときにいつでも、などに呼び出される。

【0060】

アービトレーションサイクル300は、ノード105またはノード110のいずれかが、ブロック302からブロック304に進むことによって開始される。ブロック304では、クラスタマネージャ202(図3)が、記憶サブシステム106の現在の所有権をすべて終了させる。一実施形態では、これはバス112をリセットすることによって達成される。次にこの動作によって、ストレージネットワーク100のすべての他のノードにア

アビトレーションサイクル 300 を実行させ、すべてのボリュームをオフラインモードにする。このモードでは、ボリュームマネージャ 204 が記憶サブシステム 106 のすべてのアクセスをロックする。一実施形態では、アビトレーション実行ノードは、ストレージネットワーク 100 のすべてのノードが確実にアビトレーションに入るようするために、アビトレーションサイクル 300 が進行するまでの所定の遅延期間、待機する。

【0061】

ロック 306 では、それぞれの新しいかまたは除去された記憶デバイス 106 についての構成情報を更新するために、クラスタマネージャ 202 が、ストレージネットワーク 100 内のすべての他のノードをスキャンするように、ボリュームマネージャ 204 に命令する。ロック 306 の終わりには、他のノードによって所有されていた情報が変更されている場合があるために、ボリュームマネージャ 204 によって維持される構成情報は、部分的にのみ完了することになる。したがって、ロック 308 では、クラスタマネージャ 202 が、以前にストレージネットワーク 100 のノードによって所有されていた、それらの記憶サブシステム 106 を識別するリストを生成するように、ボリュームマネージャ 204 に命令する。

【0062】

ロック 309 では、ボリュームマネージャ 204 が、生成されたリストの各記憶デバイス 106 から情報を読み取って処理する。ボリュームマネージャ 204 は、内部構成データベースを再構築する。この動作によって、アビトレーションサイクル 300 以前に、クオーラムリソース 全体が異なるノードによって所有されていた場合であっても、アビトレーション実行ノードがストレージネットワーク 100 に関するクオーラムリソース を確実に発見する。ロック 309 が終わると、ボリュームマネージャ 204 は、すべての記憶サブシステム 106 およびそのすべてのボリュームに関する情報を得る。

【0063】

次に、ロック 310 で、クラスタマネージャ 202 は、クオーラムボリューム に関する連付けられたすべての物理クオーラムリソース 111 を識別するように、ボリュームマネージャ 204 に要求する。ボリュームマネージャ 204 は、すべての記憶サブシステム 106 が物理クオーラムリソース 111 を有することを判定し、ストレージネットワーク 100 に関するクオーラムボリューム 情報を再構築する。たとえば図 1 を参照すると、確実にボリュームをオンライン上に持ち出すことができるようするために、ボリュームマネージャ 204 が記憶デバイス 108 および 109 の所有権を必要に応じて識別する。ロック 310 が完了すると、クオーラムボリューム 情報は、ストレージネットワーク 100 のすべてのノードに対して一貫したものとなる。この時点で、クラスタマネージャ 202 は記憶デバイス 108 および 109 の所有権を得ようと試みる。

【0064】

ロック 312 では、クラスタマネージャ 202 が、物理クオーラムリソース、すなわち記憶デバイス 108 および 109 をアビトレーションするために、SCSI プロトコルによって指定された技法などの、バス 112 によって提供された従来のアビトレーション技法を呼び出す。これらの従来のメカニズムが終わると、ノード 105 または 110 のいずれかが、記憶デバイス 108 および 109 の両方を所有することができるか、あるいは従来のアビトレーション技法に提示された競合条件により、物理クオーラムリソース 111 の所有権を分けることができる。

【0065】

物理クオーラムリソース 111 のアビトレーションが完了した後、ボリュームマネージャ 204 は、ローカルノード、すなわちクラスタマネージャ 202 を実行中のノードが、クオーラムボリューム に必要な記憶デバイス 108 および 109 の両方の所有権を首尾よく取得したかどうかを判定する。首尾よく取得した場合、ボリュームマネージャ 204 はクオーラムボリューム を装着し、ロック 316 でクラスタマネージャ 202 は、ローカルノードがクラスタマスターになることを宣言し、ストレージネットワーク 100 が形成

されたことを他のノードに通知する。この時点で、他のノードはアービトレーションを終了し、ストレージネットワーク100を接合する。

【0066】

ローカルノードが記憶デバイス108および109の両方の所有権を有していない場合、ボリュームマネージャ204はブロック314から318に進み、ローカルノードが任意のクオーラムボリュームリソース、すなわち記憶デバイス108または109のいずれかの所有権を取得したかどうかを判定する。ローカルノードがいずれかの所有権を有していない場合、制御はクラスタマネージャ202に渡り、これがブロック320でアービトレーションを終了して、最終的にクラスタマスターとなる他のノードからの通信を待つ。

【0067】

アービトレーション実行ノードが、記憶デバイス108および109の両方ではなく一方の所有権を有する場合、ボリュームマネージャ204はブロック318からブロック322に進み、ボリュームリストがクオーラムを形成するのに十分であるかどうかを判定する。ボリュームマネージャは、ボリュームリストが適切であるかどうかを判定する際に、単純な多数決または加重投票方式などの、いくつかの異なるアルゴリズムを使用する。ボリュームリストが不十分な場合、ボリュームマネージャ204は何らかのクオーラムリソースを解放する。クラスタマネージャ202はブロック320に進み、最終的にクラスタマスターになる他のノードからの通信を待つ。

【0068】

ただし、ボリュームマネージャ204が、ボリュームリストが十分であると判定した場合、ボリュームマネージャ204はブロック322からブロック323に進み、クオーラムボリュームの装着が安全かどうかを判定する。この判定は、ボリューム特有の情報に基づいて行われる。たとえば、クオーラムボリュームが連結されたかまたはストライピングされたエクステントを使用する場合、ボリュームマネージャ204は、常に、1つのエクステントのみが所有されているときにはクオーラムボリュームの装着は安全でないと判定する。他の例として、クオーラムボリュームがRAIDVである場合、ボリュームマネージャ204は、1つを除くすべてのエクステントが必要であるというような、「マイナス1」アルゴリズムを適用することができる。さらに、ボリュームマネージャ204はユーザが選択可能な基準を適用することもできる。たとえば、クオーラムボリュームがミラーである場合、ユーザはすべてのエクステントを要求するか、または単純な多数決を要求するように、ボリュームマネージャ204を構成することができる。ボリュームマネージャ204がクオーラムボリュームを安全に装着できる場合、ボリュームマネージャ204はクオーラムボリュームを装着し、クラスタマネージャ202はブロック316に進んで、ローカルノードをクラスタマスターとして宣言する。

【0069】

ただし、ボリュームマネージャ204がクオーラムボリュームの装着が安全に実行できないと判定した場合、クラスタマネージャ202は所定の期間待機する。ブロック326で、所定の時間内にクラスタマスターからの通信が受け取られなかった場合、クラスタマネージャ202はブロック304に戻り、本発明のアービトレーション方法を繰り返す。一実施形態では、アービトレーションサイクル300が反復されるごとに、遅延期間が増加する。

【0070】

〔結論〕

以上、ボリューム管理ならびに論理リソースの構成の基礎となる物理的特性についての知識を必要とせずに、クラスタソフトウェアが論理クオーラムリソースのアービトレーションを実行できるようにする、本発明のアービトレーション方式の様々な実施形態について説明してきた。ボリューム管理ソフトウェアは、アービトレーションプロセスを介してクラスタの所有権がどのように確立されるかを知らずに、基礎となる記憶デバイスを管理する。この方法で、本発明は、クラスタ管理の役割とボリューム管理の役割を明確に分離する。本発明は、特許請求の範囲およびその等価物によってのみ限定されるものであるこ

とが意図されている。

【図面の簡単な説明】

【図 1】

本発明の実施形態がそれと共に実施可能なハードウェアおよび動作環境を示す図である。

【図 2】

2つのコンピューティングシステムおよび様々な記憶デバイスを備えたストレージネットワークのシステムレベルの概要を示すブロック図である。

【図 3】

クラスターアービトレーションの役割とボリュームおよび基礎となる記憶デバイスの管理とを明確に分離した、協働するソフトウェア構成要素を有するソフトウェアシステムの一実施形態を示すブロック図である。

【図 4】

本発明に従ってシステムが論理クオーラムリソースをアービトレーションする、図3のソフトウェアシステムの一動作モードを示すフローチャートである。

【誤訳訂正 2】

【訂正対象書類名】図面

【訂正対象項目名】図 4

【訂正方法】変更

【訂正の内容】

【図4】

