

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2005-535961
(P2005-535961A)

(43) 公表日 平成17年11月24日(2005.11.24)

(51) Int. Cl. ⁷	F I	テーマコード (参考)
G06F 12/00	G06F 12/00 545A	5B014
G06F 3/06	G06F 12/00 514E	5B065
G06F 13/10	G06F 12/00 545B	5B082
	G06F 3/06 301A	
	G06F 13/10 340A	
審査請求 有 予備審査請求 未請求 (全 28 頁)		

(21) 出願番号 特願2004-527664 (P2004-527664)
 (86) (22) 出願日 平成15年7月28日 (2003.7.28)
 (85) 翻訳文提出日 平成17年3月15日 (2005.3.15)
 (86) 国際出願番号 PCT/US2003/023597
 (87) 国際公開番号 W02004/015521
 (87) 国際公開日 平成16年2月19日 (2004.2.19)
 (31) 優先権主張番号 10/215,917
 (32) 優先日 平成14年8月9日 (2002.8.9)
 (33) 優先権主張国 米国 (US)

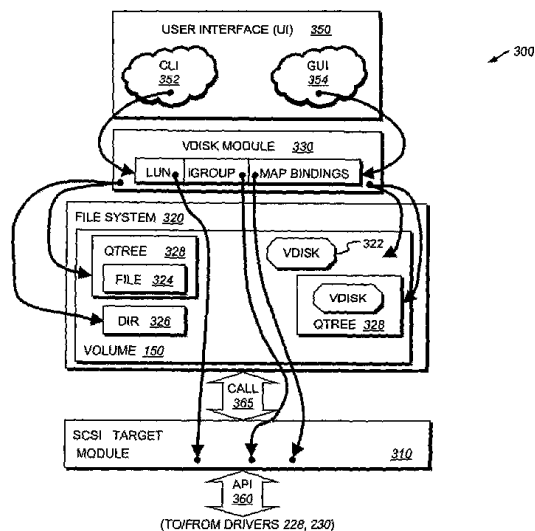
(71) 出願人 500261341
 ネットワーク・アプライアンス・インコーポレイテッド
 アメリカ合衆国94089カリフォルニア州サニーベイル、イースト・ジャーバードライブ495番
 (74) 代理人 100087642
 弁理士 古谷 聡
 (74) 代理人 100076680
 弁理士 溝部 孝彦
 (74) 代理人 100121061
 弁理士 西山 清春

最終頁に続く

(54) 【発明の名称】 ファイル・アクセス・プロトコルとブロック・アクセス・プロトコルの統合サポートを提供するマルチプロトコル・ストレージ・アプライアンス

(57) 【要約】

記憶装置に格納された情報にアクセスするためのファイル・プロトコルおよびブロック・プロトコルを、ネットワーク・アタッチド・ストレージ (NAS) デプロイメントとストレージ・エリア・ネットワーク (SAN) デプロイメントの両方に対し、統一的方法でサービス提供するマルチプロトコル・ストレージ・アプライアンス。該アプライアンスのストレージ・オペレーティング・システムは、新規の仮想化モジュールと協働するファイルシステムを実施し、記憶装置によって提供された記憶空間を「仮想化」する仮想システムを提供する。特に、ファイルシステムは、記憶装置に格納された情報をブロック・アクセスをする際に使用されるボリューム管理機能を提供する。仮想化システムは、ファイルシステムが情報を名前付きのファイル、ディレクトリおよび仮想ディスク (v d i s k) ストレージ・オブジェクトとして論理編成できるようにし、ファイルやディレクトリに対するファイル単位のアクセスを可能にするとともに、v d i s k に対するブロック単位のアクセスを更に可能にすることにより、NASアプライアンスとSANアプライ



【特許請求の範囲】

【請求項 1】

記憶装置に格納された情報にアクセスするためのファイル・プロトコルおよびブロック・プロトコルを、ネットワーク・アタッチド・ストレージ (NAS) デプロイメントと、ストレージ・エリア・ネットワーク (SAN) デプロイメントとの両方に対し、統一的態様でサービス提供するように構成されたマルチプロトコル・ストレージ・アプライアンスであって、

仮想化モジュールと協働するファイルシステムを実施し、記憶装置によって提供される記憶空間を仮想化するように構成されたストレージ・オペレーティング・システムを含む、マルチプロトコル・ストレージ・アプライアンス。

10

【請求項 2】

前記ファイルシステムは、情報をファイル、ディレクトリおよび仮想ディスク (v d i s k) として論理編成し、ファイルおよびディレクトリに対するファイル単位のアクセスを可能にするとともに、v d i s k に対するブロック単位のアクセスを可能にすることにより、NAS アプライアンスと SAN アプライアンスの統一的記憶方法を提供する、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 3】

前記仮想化モジュールは、v d i s k モジュールと、SCSI (Small Computer Systems Interface) ターゲット・モジュールとを含む、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

20

【請求項 4】

前記 v d i s k モジュールは、ファイルシステム上に層として形成され、システム管理者によってマルチプロトコル・ストレージ・アプライアンスに発行されたコマンドにตอบสนองして、管理インタフェースによるアクセスを有効にする、請求項 3 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 5】

前記管理インタフェースはユーザ・インタフェース (UI) を含む、請求項 4 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 6】

前記 v d i s k モジュールは、前記 UI を通じて発行された一連の v d i s k コマンドを実行することにより、前記 SAN デプロイメントを管理する、請求項 5 に記載のマルチプロトコル・ストレージ・アプライアンス。

30

【請求項 7】

前記 v d i s k コマンドは、原始的なファイルシステム処理に変換され、ファイルシステムおよび SCSI ターゲット・モジュールに作用し、前記 v d i s k を実施する、請求項 6 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 8】

前記 SCSI ターゲット・モジュールは、論理ユニット番号 (LUN) を v d i s k に変換するマッピング手順を提供することにより、ディスクまたは LUN のエミュレーションを開始する、請求項 7 に記載のマルチプロトコル・ストレージ・アプライアンス。

40

【請求項 9】

前記 SCSI ターゲット・モジュールは、SAN ブロック空間と、ファイルシステム空間との間の変換層を提供する、請求項 8 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 10】

前記ファイルシステムによる汎用空間管理に関し、前記仮想化記憶空間は、SAN ストレージ・オブジェクトと NAS ストレージ・オブジェクトの共存を可能にする、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 11】

前記ファイルシステムは前記仮想化モジュールと協働して仮想化システムを提供し、前

50

記仮想化記憶空間内に共存するSANストレージ・オブジェクトとNASストレージ・オブジェクトに信頼性保証を提供する、請求項10に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項12】

前記ファイルシステムは、前記記憶装置に格納された情報に対してブロック単位のアクセスをする際に使用されるボリューム管理機能を提供する、請求項1に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項13】

前記記憶装置はディスクである、請求項12に記載のマルチプロトコル・ストレージ・アプライアンス。

10

【請求項14】

前記ファイルシステムは、(i)ストレージ・オブジェクトの名前付けなどのファイルシステム・セマンティックと、(ii)ボリューム・マネージャに関連する機能とを提供する、請求項1に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項15】

前記ボリューム・マネージャに関連する機能は、

前記記憶装置を集める機能、

前記記憶装置の記憶帯域幅を集める機能、および

ミラーリングやRAID (Redundant Array of Independent Disks) などの信頼性保

証機能、

20

のうちの少なくとも1つを含む、請求項14に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項16】

マルチプロトコル・ストレージ・アプライアンスに接続された記憶装置上での情報の編成に関連するストレージ・サービスを提供するための方法であって、

マルチプロトコル・ストレージ・アプライアンス上で実行されているストレージ・オペレーティング・システムの仮想化モジュールと協働するファイルシステムを用いて、記憶装置によって提供される記憶空間を仮想化するステップと、

情報をファイル、ディレクトリおよび仮想ディスク(vdisk)として前記仮想化記憶空間内に論理編成することにより、ネットワーク・アタッチド・ストレージ(NAS)とストレージ・エリア・ネットワーク(SAN)の統一的記憶方法を提供し、前記仮想化記憶空間内の前記ファイルシステムによる汎用空間管理に関し、前記オブジェクトの共存を可能にするステップと、

30

前記マルチプロトコル・ストレージ・アプライアンスの統合型ネットワーク・プロトコル・スタックによって提供されたデータ・パスを介し、ブロック・アクセス・プロトコルおよびファイル・アクセス・プロトコルを用いて、前記記憶装置に論理編成されて格納された情報にアクセスするステップと、

からなる方法。

【請求項17】

前記仮想化記憶空間内に共存する前記ファイル、ディレクトリおよび仮想オブジェクトに対して信頼性保証を与えるステップをさらに含む、請求項16に記載の方法。

40

【請求項18】

マルチプロトコル・ストレージ・アプライアンスに接続された記憶装置上に格納される情報の編成に関連するストレージ・サービスを提供するように構成された、マルチプロトコル・ストレージ・アプライアンスのストレージ・オペレーティング・システムであって、

前記マルチプロトコル・ストレージ・アプライアンスに格納された情報をクライアントがブロック・アクセス・プロトコルおよびファイル・アクセス・プロトコルを用いてアクセスするためのデータ・パスを提供する統合型ネットワーク・プロトコル・スタックと、

50

仮想化モジュールと協働し、前記記憶装置によって提供される記憶空間を仮想化するファイルシステムと、

からなるストレージ・オペレーティング・システム。

【請求項 19】

前記ファイルシステムは、情報をファイル、ディレクトリおよび仮想ディスク (v d i s k) として論理編成することにより、ファイル単位のアクセス・プロトコルと、ブロック単位のアクセス・プロトコルとを用いた、ネットワーク・アタッチド・ストレージ (N A S) アプライアンスとストレージ・エリア・ネットワーク (S A N) の統一的方法を提供する、請求項 18 に記載のストレージ・オペレーティング・システム。

【請求項 20】

前記ブロック単位のアクセス・プロトコルは、「S C S I e n c a p u s u l a t e d o v e r T C P」(i S C S I) および「S C S I e n c a p u s u l a t e d o v e r F i b e r C h a n n e l」(F C P) のような S C S I (Small Computer Systems Interface) 系プロトコルを含む、請求項 19 に記載のストレージ・オペレーティング・システム。

10

【請求項 21】

前記統合型ネットワーク・プロトコル・スタックは、

ネットワーク・プロトコル層と、

前記ネットワーク・プロトコル層に作用し、前記ファイルシステムによって編成されたファイルおよびディレクトリに対するファイル単位のプロトコル・アクセスを提供する

20

ファイルシステム・プロトコル層と、

前記ネットワーク・プロトコル層を介して配置され、前記ファイルシステムによって編成された v d i s k に対するブロック単位のプロトコル・アクセスを提供する i S C S I ドライバと、

を含む、請求項 20 に記載のストレージ・オペレーティング・システム。

【請求項 22】

前記統合型ネットワーク・プロトコル・スタックは、前記ファイルシステム・プロトコル層のファイル・アクセス・プロトコルに対する直接アクセス搬送機能を提供する仮想インタフェース層をさらに含む、請求項 21 に記載のストレージ・オペレーティング・システム。

30

【請求項 23】

前記統合型ネットワーク・プロトコル・スタックは、前記ファイルシステムによって編成された v d i s k にアクセスするためのブロックアクセス要求を送受信するように構成されたファイバ・チャネル (F C) ドライバをさらに含む、請求項 21 に記載のストレージ・オペレーティング・システム。

【請求項 24】

前記 F C ドライバおよび i S C S I ドライバは、前記 v d i s k に対して F C 特有のアクセス制御および i S C S I 特有のアクセス制御を提供し、マルチプロトコル・ストレージ・アプライアンス上の v d i s k にアクセスするときに、i S C S I および F C P に対する v d i s k のエクスポートの管理を提供する、請求項 23 に記載のストレージ・オペレーティング・システム。

40

【請求項 25】

マルチプロトコル・ストレージ・アプライアンスの記憶装置に格納された情報にアクセスするためのファイル・プロトコルおよびブロック・プロトコルを、ネットワーク・アタッチド・ストレージ (N A S) デプロイメントと、ストレージ・エリア・ネットワーク (S A N) デプロイメントとの両方に対し、統一の態様でサービス提供する方法であって、

(i) 前記アプライアンスを第 1 のネットワークに接続するネットワーク・アダプタと

、(ii) N A S クライアントがファイルとして格納された情報にアクセスするために発行したファイル単位の要求に対し、前記アプライアンスが応答できるようにするためのファイルシステム機能とを用いて、N A S サービスを提供するステップと、

50

(i)前記アプライアンスを第2のネットワークに接続するネットワーク・アダプタと、SANクライアントが仮想ディスク(vdisk)として格納された情報にアクセスするために発行したブロックベースの要求に対し、前記アプライアンスが応答できるようにするためのボリューム管理機能とを用いて、SANサービスを提供するステップと、
からなる方法。

【請求項26】

前記マルチプロトコル・ストレージ・アプライアンスに格納されたファイルおよびvdiskの名前による管理を提供することにより、ファイル単位およびブロック単位のストレージに対して一様な名前付けを提供するステップと、

前記記憶装置に格納された名前付きファイルおよびvdiskの階層構造を提供するステップと、

【請求項27】

マルチプロトコル・ストレージ・アプライアンスの記憶装置に格納されたストレージ・オブジェクトにアクセスするためのファイル・プロトコルおよびブロック・プロトコルを、ネットワーク・アタッチド・ストレージ(NAS)デプロイメントと、ストレージ・エリア・ネットワーク(SAN)デプロイメントとの両方に対し、統一的態様でサービス提供する方法であって、

前記記憶装置を汎用記憶空間内を表す1以上のボリュームとして編成するステップと、

前記汎用記憶空間内でのSANストレージ・オブジェクトとNASストレージ・オブジェクトの共存を可能にするステップと、

前記ストレージ・アプライアンスのマルチプロトコル・エンジンにおいて、前記SANストレージ・オブジェクトおよび前記NASストレージ・オブジェクトにアクセスするためのブロック単位の要求およびファイル単位の要求を受信するステップと、

前記ブロック単位の要求およびファイル単位の要求に応答して、前記SANストレージ・オブジェクトおよびNASストレージ・オブジェクトにアクセスし、それらを返すステップと、

からなる方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はストレージ・システムに関し、詳しくは、ファイル・プロトコルとブロック・プロトコルをサポートするマルチプロトコル・ストレージアプライアンスに関する。

【背景技術】

【0002】

ストレージ・システムとは、メモリ、テープ、ディスクのような書き込み可能な持続性記憶装置における情報編成に関するストレージサービスを提供するコンピュータである。ストレージ・システムは一般に、ストレージ・エリア・ネットワーク(SAN)環境や、ネットワーク・アタッチド・ストレージ(NAS)環境に配置される。NAS環境で使用される場合、ストレージ・システムはオペレーティングシステムを含むファイルサーバとして実施され、情報をディレクトリやファイルの階層構造として(例えばディスク上に)論理編成するファイルシステムを実施する。「ディスク上」の各ファイルは、ファイルの実際のデータ等の情報を格納するように構成された、例えばディスクブロックのような一連のデータ構造として実施される。一方、ディレクトリは、他のファイルやディレクトリに関する情報を格納する特別な形のファイルとして実施される。

【0003】

さらに、ファイルサーバまたはファイラは、クライアント/サーバモデルの情報配送に従って動作するように構成される。そのため、多数のクライアントシステム(クライアント)が、ファイラに格納されたファイル等の共有リソースにアクセスすることができる。ファイルの共有はNASシステムの特徴であり、これはファイルやファイルシステムに対

10

20

30

40

50

する意味レベルのアクセスによって実施される。N A Sシステム上の情報ストレージは一般に、地理的に分散されたイーサネット（登録商標）のような相互接続通信リンクの集まりを含むコンピュータネットワーク全体にわたって配置される。クライアントは一般に、T C P / I Pのような所定のプロトコルに従って個々のデータフレームやデータパケットをやり取りすることによってファイラと通信する。

【 0 0 0 4 】

クライアント/サーバモデルの場合、クライアントは、ポイント・ツー・ポイントリンク、共有ローカルエリアネットワーク、ワイドエリアネットワーク、あるいはインターネット等の公共ネットワーク上に実施された仮想私設ネットワークなどのコンピュータネットワークを介してクライアントをファイラに「接続」するためのアプリケーションを、コンピュータ上で実行している。N A Sシステムは一般に、ファイル単位のアクセス・プロトコルを使用している。従って各クライアントは、アクセスしたい1以上のファイルをそのデータがディスク上に格納されている位置（例えばブロック）とは無関係に識別し、ネットワークを介してファイルシステム・プロトコル・メッセージを（パケットの形で）ファイルシステムに発行することにより、ファイラにサービスを要求する。従来のC I F S (Common Internet File System) プロトコル、N F S (Network File System) プロトコル、D A F S (Direct Access File System) プロトコルといった複数のファイルシステムプロトコルをサポートするために、クライアントと通信するファイラの機能は拡張される場合もある。

10

【 0 0 0 5 】

S A Nは、ストレージ・システムと該ストレージ・システムが有する記憶装置との間にダイレクト接続を確立することが可能な高速ネットワークである。S A Nは、ストレージ・バスに対する拡張のように見える。そのため、ストレージ・システムのオペレーティングシステムは、「拡張バス」を介して、ブロック単位のアクセス・プロトコルにより、格納された情報にアクセスすることができる。その場合、拡張バスは一般に、「S C S I e n c a p s u l a t e d o v e r F C」または「S C S I e n c a p s u l a t e d o v e r T C P / I P / E t h e r n e t」のようなブロックアクセスプロトコルを用いて動作するように構成されたファイバチャネル（F C）またはイーサネット（登録商標）媒体（すなわちネットワーク）として実施される。

20

【 0 0 0 6 】

S A N構成すなわちS A Nデプロイメントでは、アプリケーションサーバ等のストレージ・システムからストレージを切り離し、ある程度の情報ストレージをアプリケーションサーバレベルで共有することができる。しかしながら、S A Nが1つのサーバに対して専用になっている環境もある。S A N環境の中には、情報をデータベース形態に編成するものもあれば、ファイル単位の編成を使用するものもある。情報がファイルとして編成される場合、情報を要求するクライアントは、ファイルマッピングを保持し、ファイルセマンティックを管理するとともに、その要求（およびサーバ応答）をディスク上のブロックアドレスによって、例えば論理ユニット番号（L U N）によってアドレス指定する。

30

【 0 0 0 7 】

従来の方法は一般に、2つの異なる手法を用いてS A N環境とN A S環境に対処する。両方の環境に対して1つの手法を提供するそれらの方法の場合、N A S機能は通常、S A Nストレージ・プラットフォームに取り付けられた例えば「サイドカー」デバイスを用いて、S A Nストレージ・プラットフォーム上に配置される。しかしながら、それらの従来システムは通常、ストレージをS A Nストレージ・ドメインとN A Sストレージ・ドメインのそれぞれに分割してしまう。すなわち、S A Nドメインの記憶空間とN A Sドメインの記憶空間が共存することはなく、それらが、例えばユーザ（システム管理者）によって実施される構成プロセスにより物理的に区分される。

40

【 0 0 0 8 】

そのような従来システムの一例は、E M C（登録商標）社から市販されているS y m m e t r i x（登録商標）s y s t e mプラットフォームである。大まかに言えば、S A

50

Nストレージ・システム (Symmetrix system) の個々のディスクは、N A S ・ サイドカー ・ デバイス (例えば、Celerra (登録商標) デバイス) に割り当てられ、次いでそれらのディスクが、例えば NFS プロトコルや CIFS プロトコルによって N A S クライアントにエクスポートされる。システム管理者は、「ユーザ定義ボリューム」を構築するために集めるそれらのディスクの「スライス」(エクステント)の位置およびディスク数を決定した後、それらのボリュームをどう使うかを決定する。従来から S A N 環境で使用されている「ボリューム」という用語は、物理ディスクと、それらのディスク内のエクステントとを指定し、それらのエクステント/ディスクをユーザ定義ボリュームエンティティに結合させる処理によって構築されたストレージ・エンティティを意味する。注意すべき点は、ユーザ定義ボリュームを含む S A N 系ディスクと N A S 系ディスクとが、システムプラットフォーム内で物理的に区分される点である。

10

【0009】

一般にシステム管理者は、その決定を、システムの基礎物理態様に関する知識を有するユーザ向けの複雑なユーザ・インタフェースを通じて実施する。つまり、クライアントに代わって S A N プラットフォームのビューを提供するために、システム管理者がユーザ・インタフェースで操作しなければならないものは、物理的ディスク構造や管理が中心となる。例えば、ユーザ・インタフェースは、物理ディスクの指定を促すとともに、ユーザ定義ボリュームを構築するのに必要とされるそれらのディスク内のエクステントのサイズの指定を促す場合がある。さらに、ユーザ・インタフェースは、そうしたエクステントやディスクの物理的位置の指定を促すだけでなく、それらを「相互接合」(編成)し、それらをディスクまたは L U N に対応するユーザ定義ボリュームとして S A N クライアントから見えるようにする(エクスポートする)態様の指定を促す。物理ディスクおよびそれらのエクステントが選択され、ボリュームが構築された後、それらのディスク/エクステントだけが、そのボリュームを含む。システム管理者は、構築されたボリュームに対し、信頼性の形、例えば、R A I D (Redundant Array of Independent Disks) 保護レベルおよび/またはミラーリングなどをさらに指定しなければならない。そして、R A I D グループが、それらの選択されたディスク/エクステントの最上部に置かれる。

20

【0010】

要するに、従来システムの方法では、システム管理者が、ディスクの物理レイアウトやディスクの編成を詳細に設定し、S A N クライアントに対して1つの L U N としてエクスポートされるユーザ定義ボリュームを作成しなければならない。この従来の方法に関連する管理はすべて、物理ディスク単位で行なわれる。システム管理者がユーザ定義ボリュームのサイズを拡大する場合、ディスクを追加し、ボリュームを構成するディスクに格納されたデータに関連する冗長情報を含めるように R A I D 計算を再計算する。明らかに、これは複雑でコストのかかる方法である。

30

【発明の開示】**【発明が解決しようとする課題】****【0011】**

本発明は、S A N ストレージ環境および N A S ストレージ環境に対し、簡単で効率的な統合型ソリューションを提供することを目的としている。

40

【課題を解決するための手段】**【0012】**

本発明は、記憶装置に格納された情報にアクセスするためのファイル・プロトコルおよびブロック・プロトコルを、ネットワーク・アタッチド・ストレージ (N A S) 環境とストレージ・エリア・ネットワーク (N A S) 環境の両方に対し、統一された態様でサービス提供するマルチプロトコル・ストレージ・アプライアンスに関する。このアプライアンスのストレージ・オペレーティングシステムは、新規の仮想化モジュールと協働するファイルシステムを実施し、記憶装置によって提供される記憶空間を「仮想化」する仮想化システムを提供する。特に、ファイルシステムは、記憶装置に格納された情報に対するブロック単位のアクセスに使用されるボリューム管理機能を提供する。仮想化システムは、フ

50

ファイルシステムが、情報を名前付きのファイル、ディレクトリおよび仮想ディスク (v d i s k) ・ストレージ・オブジェクトとして論理編成できるようにし、ファイルやディレクトリに対するファイル単位のアクセスを可能にすると同時に、v d i s kに対するブロック単位のアクセスをさらに可能にすることにより、N A SアプライアンスとS A Nアプライアンスの統一的記憶方法を提供する。

【0013】

実施例において、仮想化モジュールは、例えばv d i s kモジュールおよびS C S I (Small Computer Systems Interface)ターゲット・モジュールとして実施される。v d i s kモジュールは、ブロックベースのS C S Iターゲット・モジュールからファイルシステムによって管理されるブロックへのデータ・パスを提供する。また、v d i s kモジュールは、システム管理者によってマルチプロトコル・ストレージ・アプライアンスに発行されたコマンドに回答し、ファイルシステムに作用し、簡易ユーザ・インタフェース (U I) などの管理インタフェースによるアクセスを有効にする。さらに、v d i s k管理モジュールは、とりわけ、システム管理者によってU Iを通じて発行された一連のv d i s kコマンドを実行することにより、S A Nデプロイメントを管理する。それらのv d i s kコマンドは原始的なファイルシステム処理に変換され、ファイルシステムおよびS C S Iターゲット・モジュールに作用し、v d i s kを実施する。

10

【0014】

次にS C S Iターゲット・モジュールは、アクセス要求に対する応答として、アクセス要求の中で指定されたL U Nに対する論理ブロックアクセスをv d i s kに対する仮想ブロックアクセスに変換し、v d i s kをL U Nに変換する。次にS C S Iターゲット・モジュールは、S A Nブロック (L U N) 空間と、L U Nがブロックとして表現されるファイルシステム空間との間の、仮想化システムの変換層として機能する。S A N仮想化をファイルシステム上に「配置」することにより、マルチプロトコル・ストレージ・アプライアンスは、従来のシステムで行なわれている方法を逆にすることにより、実質的に全てのストレージ・アクセス・プロトコルについて単一の統一されたストレージプラットフォームを提供する。

20

【0015】

有利なことに、統合型マルチプロトコル・ストレージ・アプライアンスは、データの完全性を維持しつつ、アクセス制御を提供し、適当であれば、全てのプロトコルについてファイルおよびv d i s kの共有を提供する。さらに、ストレージ・アプライアンスは、N A Sストレージ・オブジェクトやS A Nストレージ・オブジェクトを作成するときにユーザがストレージリソースを配分する必要性を無くす埋め込み型/一体型仮想化機能を提供する。それらの機能は、アプライアンス内の汎用空間管理に関し、S A NオブジェクトとN A Sオブジェクトとを共存させることが可能な仮想化記憶空間を含む。さらに、統合型マルチプロトコル・ストレージ・アプライアンスは、同一ディスクに対する複数のブロックアクセスプロトコルの同時サポートを提供するだけでなく、クラスタリングをサポートする異種のS A N環境も提供する。

30

【0016】

本発明の上記の利点およびその他の利点は、添付の図面とともに下記の説明を参照することで、さらによく理解できるであろう。図中、同一の参照符号は、同一要素または機能的に類似の要素を指している。

40

【発明を実施するための最良の形態】**【0017】**

本発明は、記憶装置に格納された情報にアクセスするためのファイル・プロトコルおよびブロックプロトコルを、統一的態様でサービス提供するマルチプロトコル・ストレージ・アプライアンスに関する。これに関し、統合型マルチプロトコル・アプライアンスとは、ネットワーク・アタッチド・ストレージ (N A S) デプロイメントやストレージ・エリア・ネットワーク (S A N) のユーザ (システム管理者) およびクライアントにとって記憶空間が再利用可能であることを含めて、ストレージ・サービスの管理が容易であること

50

や、記憶装置の再構成が簡単であること等の特徴を有するコンピュータを意味する。ストレージ・アプライアンスは、ファイルシステムを通してNASサービスを提供するとともに、SAN仮想化を通じて論理ユニット番号(LUN)のエミュレーションを含むSANサービスを提供する。

【0018】

図1は、ディスク130などの記憶装置上での情報編成に関するストレージ・サービスを提供するように構成されたマルチプロトコル・ストレージ・アプライアンス100を示す略ブロック図である。ストレージ・アプライアンス100は例えば、システム・バス123によって相互接続された、プロセッサ122、メモリ124、複数のネットワーク・アダプタ125, 126、ストレージ・アダプタ128を含むストレージ・システムとして実施される。また、マルチプロトコル・ストレージ・アプライアンス100は、情報を名前付きのディレクトリ、ファイルおよび仮想ディスク(vdisk)・ストレージ・オブジェクトの階層構造としてディスク130上に論理編成する仮想化システム(および、特にファイルシステム)として機能するストレージ・オペレーティング・システム200をさらにも含む。

10

【0019】

NAS系ネットワーク環境のクライアントは、ファイルの視点からストレージを見るのに対し、SAN系ネットワーク環境のクライアントは、ブロックまたはディスクの視点からストレージを見る。そのため、マルチプロトコル・ストレージ・アプライアンス100は、LUNまたはvdiskオブジェクトを作成することにより、ディスクをSANクライアントに見せる(エクスポートする)。vdiskオブジェクト(以下「vdisk」)は、仮想化システムによって実施される特殊なファイルタイプであり、SANクライアントから見るとエミュレート・ディスクとして解釈される。後で詳しく説明するように、マルチプロトコル・ストレージ・アプライアンスは、エクスポートを制御することにより、SANクライアントがそれらのエミュレート・ディスクにアクセスできるようにする。

20

【0020】

図示の実施形態において、メモリ124は、本発明に係るソフトウェアプログラムやデータ構造を格納するための、プロセッサやアダプタによってアドレス指定可能な複数の格納場所を有する。そして、プロセッサおよびアダプタは、ソフトウェアプログラムを実行し、データ構造を操作するように構成された処理要素および/または論理回路を含む。ストレージ・オペレーティング・システム200の一部は通常、メモリに常駐し、それらの処理要素によって実行され、特に、ストレージ・アプライアンスによって実施されるストレージサービスを支える記憶処理を実施することにより、ストレージ・アプライアンスの機能を構成する。本明細書に記載されている本発明のシステムおよび方法に係るプログラム命令を記憶または実行するために、他の処理手段や、種々のコンピュータ読取可能媒体のような他の記憶手段を使用してもよいことは、当業者にとって明らかであろう。

30

【0021】

ネットワークアダプタ125は、以下で例示的なイーサネット(登録商標)ネットワーク165と呼ぶポイント・ツー・ポイントリンク、ワイドエリアネットワーク、公共ネットワーク上で実施される仮想私設ネットワーク、または共有ローカルエリアネットワークを介して、ストレージ・アプライアンスを複数のクライアント160a, 160bに接続する。従って、ネットワークアダプタ125は、ストレージ・アプライアンスを従来のイーサネット(登録商標)・スイッチ170のようなネットワークスイッチに接続するために必要となる機械的電気的信号回路を備えたネットワークインタフェースカード(NIC)を含む。このNAS系ネットワーク環境の場合、クライアントは、マルチプロトコル・ストレージ・アプライアンス上にファイルとして格納された情報にアクセスするように構成される。クライアント160は、TCP/IPのような所定のプロトコルに従って個々のデータフレームまたはデータパケットをやり取りすることにより、ネットワーク165を介してストレージ・アプライアンスと通信する。

40

50

【0022】

クライアント160は、UNIX（登録商標）やMicrosoft Windows（登録商標）のような種々のオペレーティング・システム上でアプリケーションを実行するように構成された汎用コンピュータであってよい。クライアント・システムは一般に、NAS系ネットワークを介して（ファイルやディレクトリの形の）情報にアクセスする場合、ファイル単位のアクセス・プロトコルを使用する。従って、各クライアント160は、ネットワーク165を介してファイル・アクセス・プロトコル・メッセージを（パケットの形で）ストレージ・アプライアンス100に発行することにより、ストレージ・アプライアンスにサービスを要求する。例えば、Windows（登録商標）オペレーティング・システムを実行しているクライアント160aは、「CIFS over TCP/IP」プロトコルを用いてストレージ・アプライアンス100と通信することができる。一方、UNIXオペレーティング・システムを実行しているクライアント160bは、「NFS over TCP/IP」プロトコルを用いて、若しくは「RDMA over VI」プロトコルに従って転送を行なう「DAFS over VI」プロトコルを用いて、マルチプロトコル・ストレージ・アプライアンスと通信することができる。他のタイプのオペレーティング・システムを実行する他のクライアントが、他のファイル・アクセス・プロトコルを用いて統合型マルチプロトコル・ストレージ・アプライアンスと通信することも可能であることは、当業者にとって明らかである。

10

【0023】

ストレージ・ネットワーク・「ターゲット」・アダプタ126も、マルチプロトコル・ストレージ・アプライアンス100をクライアント160に接続する。クライアント160はさらに、ブロックまたはディスクとして格納された情報にアクセスするように構成される場合もある。SAN系ネットワーク環境の場合、ストレージ・アプライアンスは、ファイバ・チャネル（FC）・ネットワーク185に接続される。FCとは、SANデプロイメントで主に使用される適当なプロトコルや媒体が記載されているネットワーク規格である。ネットワーク・ターゲット・アダプタ126は、ストレージ・アプライアンス100を従来のFCスイッチ180のようなSAN・ネットワーク・スイッチに接続するために必要となる機械的電気的信号回路を備えたFCホストバスアダプタ（HBA）を含む。FC HBAは、FCアクセスを可能にするだけでなく、ストレージ・アプライアンスのファイバ・チャネル・ネットワーク処理動作の負荷も軽減する。

20

30

【0024】

クライアント160は一般に、SAN系ネットワークを介して情報にアクセスするときSCSIプロトコルのようなブロック単位のアクセス・プロトコルを使用する。SCSIは、装置に依存しない標準的プロトコルを備えた周辺装置入出力（I/O）インタフェースであり、ディスク130のような種々の周辺装置をストレージアプライアンス100に接続できるようにする。SAN環境で動作するクライアント160は、SCSI用語で「イニシエータ」と呼ばれ、データの要求やコマンドを開始する。従って、マルチプロトコル・ストレージ・アプライアンスは、要求/応答プロトコルに従ってイニシエータによって発行された要求に回答するように構成された「ターゲット」である。FCプロトコルによれば、イニシエータおよびターゲットは、終点アドレスを有し、ワールド・ワイド・

40

【0025】

マルチプロトコル・ストレージ・アプライアンス100は、「SCSI encapsulated over TCP (iSCSI)」や、「SCSI encapsulated over FC (FCP)」のような、SAN環境で使用される種々のSCSI系のプロトコルをサポートする。従って、イニシエータ（以下、クライアント160）は、ネットワーク185を介してiSCSIメッセージやFCPメッセージを発行することにより、ターゲット（以下、ストレージ・アプライアンス100）にサービスを要求し、ディスク上に格納された情報にアクセスする。クライアントが他のブロック・アクセス・

50

プロトコルを用いて一体化マルチプロトコル・ストレージ・アプライアンスにサービスを要求してもよいことは、当業者にとって明らかであろう。マルチプロトコル・ストレージ・アプライアンスは、複数のブロック・アクセス・プロトコルをサポートすることで、異種のSAN環境においても、v d i s k / L U Nに対する統一的で一貫性のあるアクセス方法を提供することができる。

【0026】

ストレージ・アダプタ128は、ストレージ・アプライアンス上で実行されているストレージ・オペレーティング・システム200と協働し、クライアントから要求された情報を取得する。この情報は、ディスク130に格納されている場合もあれば、情報を格納するように構成された他の似たような媒体に格納されている場合もある。ストレージ・アダプタは、従来の高性能FCシリアルリンクトポロジのようなI/O相互接続構成を介してストレージ・アダプタをディスクに接続するためのI/Oインタフェースを有する。情報は、ストレージ・アダプタによって取得された後、システム・バス123を介してネットワーク・アダプタ125、126に転送される前に、必要に応じてプロセッサ122（若しくはアダプタ128自体）によって処理される。そしてネットワーク・アダプタ125、126は、その情報をパケットまたはメッセージの形にフォーマットし、クライアントに返送する。

10

【0027】

アプライアンス100に対する情報の格納は、ディスク空間の全体的論理構成を定義しつつ、物理ストレージディスク130の集まりからなる1以上のストレージボリューム（例えば、VOL1~2 150）として実施することが好ましい。ボリューム内のディスクは通常、1以上のRAIDグループに編成される。RAID実施形態は、RAIDグループ内の所与の数の物理ディスクにわたってデータを「ストライプ」状に書き込み、さらにそのストライプ状のデータに関する冗長情報を適当に記憶することによって、データ記憶の信頼性/完全性を向上させる。記憶装置が故障した場合でも、冗長情報によって、失われたデータを復元することが可能になる。本発明に従ってミラーリング等の他の冗長技術を使用することも可能であることは、当業者にとって明らかであろう。

20

【0028】

具体的には、各ボリューム150は、RAIDグループ140、142、144として編成された物理ディスク130のレイから構成される。RAID4レベル構成の場合、各RAIDグループの物理ディスクは、ストライプ・データ(D)を格納するように構成された物理ディスクと、それらのデータのパリティ(P)を格納するように構成された物理ディスクとを含む。ただし、他のRAIDレベル構成（例えば、RAID5など）も可能であるものと考えられる。図示の実施形態では、1台のパリティディスクと、1台のデータディスクとからなる最小構成を採用する場合がある。しかしながら、典型的な実施形態は、1つのRAIDグループ当たり3台のデータディスクと1台のパリティディスクとを含み、1つのボリューム当たり少なくとも1つのRAIDグループを含む実施形態であろう。

30

【0029】

ディスク130に対するアクセスを容易にするために、ストレージ・オペレーティング・システム200は、新規の仮想化システムのWrite-anywareファイルシステムを実施し、ディスク130によって提供される記憶空間を「仮想化」する機能を提供する。ファイルシステムは、情報を名前付きのディレクトリおよびファイル・オブジェクト（以下、「ディレクトリ」および「ファイル」）の階層構造としてディスク上に編成する。「ディスク上」の各ファイルは、データ等の情報を格納するように構成された一連のディスク・ブロックとして実施される一方、ディレクトリは、他のファイルまたはディレクトリの名称や、それらに対するリンクを格納する特別な形のファイルとして実施される。仮想化システムは、ファイルシステムが、情報を名前付きv d i s kの階層構造としてディスク上にさらに論理編成できるようにし、名前付きのファイルおよびディレクトリに対するファイル単位のアクセス(NAS)を可能にすると同時に、ファイルベースのスト

40

50

レージ・プラットフォーム上の `v d i s k` に対するブロック単位のアクセス (`S A N`) をさらに可能にすることにより、`N A S` アプライアンスと `S A N` アプライアンスの統一的記憶方法を提供する。ファイルシステムは、`S A N` デプロイメントにおける基礎物理ストレージの構成の複雑さを単純化する。

【 0 0 3 0 】

上記のように、`v d i s k` は、普通 (通常) のファイルのうち、ディスクのエミュレーションを支援するエクスポート制御および動作制限に関連するファイルに由来する、ボリューム内の特別なファイルタイプである。`v d i s k` は、クライアントにより例えば `N F S` プロトコルや `C I F S` プロトコルを用いて作成される可能性があるファイルとは異なり、例えばユーザ・インタフェース (`U I`) を介して特別なタイプのファイル (オブジェクトとしてマルチプロトコル・ストレージ・アプライアンス上に作成される。例えば、`v d i s k` は、データを保持する特別なファイル `i n o d e` と、セキュリティ情報のような属性を保持する少なくとも1つの関連ストリーム `i n o d e` とを含む、マルチ `i n o d e` オブジェクトである。特別なファイル `i n o d e` は、エミュレート・ディスクに関連するアプリケーション・データのようなデータを格納するためのメインコンテナとして機能する。ストリーム `i n o d e` には、例えば何回ものリブート処理の間、`L U N` をエクスポートし続けるとともに、`v d i s k` を `S A N` クライアントに関連する1つのディスク・オブジェクトとして管理できるようにする属性が格納される。本発明に使用するのに都合がよい `v d i s k` およびそれに関連する `i n o d e` の一例については、「Storage Virtualization by Layering Vdisks on a File System」と題する本発明と同じ譲受人の同時係属中の米国特許出願第 (1 1 2 0 5 6 - 0 0 6 9) 号に記載されており、この特許出願は、ここで参照することにより完全に説明されたものとして本明細書に取り込まれる。

10

20

【 0 0 3 1 】

図示の実施形態において、ストレージ・オペレーティング・システムは、カリフォルニア州、サニーベイルにあるネットワーク・アプライアンス・インコーポレイテッドから市販されている `N e t A p p D a t a O N T A P` (登録商標) オペレーティング・システムであることが好ましい。`N e t A p p D a t a O N T A P` (登録商標) オペレーティング・システムは、`W A F L` (登録商標) (Write Anywhere File Layout) ファイルシステムを実施する。ただし、書き込み場所固定式ファイルシステムのような、何らかの適当なオペレーティング・システムを、本発明に記載された発明の原理に従って拡張して使用することも可能である。従って、「`W A F L`」という用語が使用された場合、その用語は、本発明の教示に適合する任意のストレージ・オペレーティング・システムを指すものとして広い意味で解釈しなければならない。

30

【 0 0 3 2 】

「ストレージ・オペレーティング・システム」という用語が本明細書で使用される場合、その用語は通常、データアクセスの管理を行なう、コンピュータ上で動作するコンピュータ実行可能コードを意味し、マルチプロトコル・ストレージ・アプライアンスの場合、ストレージ・オペレーティング・システムは、マイクロカーネルとして実施される `D a t a O N T A P` ストレージ・オペレーティング・システムのように、データ・アクセス・セマンティックを実施する場合がある。また、ストレージ・オペレーティング・システムは、`U N I X` (登録商標) や `W i n d o w s` (登録商標) などの汎用オペレーティング・システム上で動作するアプリケーションプログラムとして実施することもでき、あるいは、本明細書に記載されるようなストレージ・アプリケーションに合わせて構成される機能構成可能な汎用オペレーティング・システムとして実施することもできる。

40

【 0 0 3 3 】

さらに、本明細書に記載される本発明のシステムおよび方法が、スタンドアロンのコンピュータやその一部を含むストレージ・システムとして実施される、またはストレージ・システムを含む、いかなるタイプの特殊用途のコンピュータ (例えば、ストレージ・サービス・アプライアンスなど) および汎用コンピュータにも適用できるものであることは、当業者にとって明らかであろう。さらに、本発明の教示は、限定はしないが、`N A S` 環境

50

、SAN環境、および、クライアントまたはホストコンピュータに直接取り付けられたディスク・アセンブリなどを含む、種々のストレージ・システム・アーキテクチャに適合させることができる。従って、「ストレージ・システム」という用語は、他の装置またはシステムに接続され、ストレージ機能を実施するように構成された何らかのサブシステムだけでなく、それらの構成も含むものとして広い意味で解釈しなければならない。

【0034】

図2は、本発明で使用するのに都合がよいストレージ・オペレーティング・システム200を示す略ブロック図である。ストレージ・オペレーティング・システムは、統合型ネットワーク・プロトコル・スタックを形成するように構成された一連のソフトウェア層を含む。すなわち、一般には、ストレージ・オペレーティング・システムは、マルチプロトコル・ストレージ・アプライアンス上に格納された情報をクライアントがブロック・アクセス・プロトコルやファイル・アクセス・プロトコルを用いてアクセスできるようにするためのデータパスを提供するマルチプロトコル・エンジンを含む。プロトコルスタックは、ネットワーク・ドライバ（例えば、ギガビット・イーサネット（登録商標）・ドライバなど）からなるメディア・アクセス層210を含み、メディア・アクセス層210は、IP層212や、その支援搬送手段であるTCP層214およびユーザ・データグラム・プロトコル（UDP）層216などのネットワーク・プロトコル層と連絡する。マルチプロトコル・ファイル・アクセス機能を提供するファイルシステム・プロトコル層は、その目的のために、DAFSプロトコル218、NFSプロトコル220、CIFSプロトコル222およびハイパーテキスト・トランスファ・プロトコル（HTTP）プロトコル224を含む。VI層226は、VIアーキテクチャを実施し、DAFSプロトコル218に必要とされるRDMAのようなDAT（Direct Access Transport）機能を提供する。

【0035】

iSCSIドライバ層228は、TCP/IPネットワークプロトコル層を介したブロック・プロトコル・アクセス機能を提供し、FCドライバ層230は、FC HBA126とともに動作し、ブロックアクセス要求の送受信と、統合型ストレージ・アプライアンスに対する応答とを行なう。FCドライバおよびiSCSIドライバは、LUN（vdisk）に対して、FC特有のアクセス制御およびiSCSI特有のアクセス制御を行なう。すなわち、FCドライバおよびiSCSIドライバは、マルチプロトコル・ストレージ・アプライアンス上の1つのvdiskをアクセスするとき、iSCSIとFCPのどちらか一方に対するvdiskのエクスポート、または、両方に対するvdiskのエクスポートを管理する。さらに、ストレージ・オペレーティング・システムは、RAIDプロトコルなどのディスク・ストレージ・プロトコルを実施するディスク・ストレージ層240と、SCSIプロトコルなどのディスク・アクセス・プロトコルを実施するディスク・ドライバ層250とを含む。

【0036】

それらのディスク・ソフトウェア層を統合型ネットワーク・プロトコル・スタック層に橋渡しするのが、本発明による仮想化システム300である。図3は、仮想化モジュールと通信するファイルシステム320によって実施される仮想化システム300を示す略ブロック図である。仮想化モジュールは、例えばvdiskモジュール330およびSCSIターゲット・モジュール310として実施される。vdiskモジュール330、ファイルシステム320およびSCSIターゲット・モジュール310が、ソフトウェアで実施することも、ハードウェアで実施することも、ファームウェアで実施することもでき、さらにそれらの組み合わせで実施することも可能であることは、当業者にとって明らかであろう。vdiskモジュール330は、ファイルシステム320上に階層化され（およびファイルシステム320と通信し）、ブロックベースのSCSIターゲット・モジュールからファイルシステムによって管理されるブロックへのデータ・パスを提供する。また、vdiskモジュールは、システム管理者がマルチプロトコル・ストレージ・アプライアンスに対して発行したコマンドに回答して、簡易ユーザ・インタフェース（UI350）のような管理インタフェースによるアクセスを有効にする。要するに、vdiskモジ

10

20

30

40

50

ユーザ 330 は、システム管理者が UI 350 を通じて発行した様々な一連の `vdisk` コマンドを実施することにより、特に SAN デプロイメントを管理する。それらの `vdisk` コマンドは、原始的なファイルシステム処理（「プリミティブ」）に変換され、ファイルシステム 320 および SCSI ターゲット・モジュール 310 に作用し、`vdisk` を実施する。

【0037】

次いで SCSI ターゲット・モジュール 310 は、アクセス要求に対し、そのアクセス要求の中で指定された LUN に対する論理ブロックアクセスを特別な `vdisk` ファイルタイプに対する仮想ブロックアクセスに変換し、`vdisk` を LUN に変換する。SCSI ターゲット・モジュールは、例えば FC ドライバや iSCSI ドライバ 228 と、ファイルシステム 320 との間に配置され、仮想化システム 300 にとって、SAN ブロック（LUN）空間とファイルシステム空間との間の変換層として機能する。その際、LUN は `vdisk` 322 として表現される。マルチプロトコル・ストレージ・アプライアンスは、SAN 仮想化をファイルシステム 320 全体に「配置」することにより、従来のシステムで行なわれていた方法を逆転させ、実質的に全てのストレージ・アクセス・プロトコルに対して 1 つの統一されたストレージ・プラットフォームを提供する。

【0038】

本発明によれば、ファイルシステムは、ディスクなどの記憶装置に格納された情報に対してファイル単位のアクセスをする際に使用される機能を提供する。さらに、ファイルシステムは、格納された情報に対してブロック単位のアクセスをする際に使用されるボリューム管理機能を提供する。すなわち、ファイルシステム 320 は、ファイルシステム・セマンティック（ストレージを個々のオブジェクトに区別することや、それらのストレージ・オブジェクトの名前付けなど）を提供する他に、通常ならばボリューム・マネージャに関連する機能も提供する。そのような機能には、(i) ディスクを集める機能、(ii) ディスクの記憶帯域幅を集める機能、(iii) ミラーリングおよび/またはパリティ（RAID）のような信頼性保証機能などがあり、それらによって 1 以上のストレージ・オブジェクトをファイルシステム上に階層化して表現する。マルチプロトコル・ストレージ・アプライアンスの特徴は、マルチプロトコル・ストレージ・アプライアンスを特に SAN デプロイメントで使用したときの、それらのボリューム管理機能に関する使用の簡単さにある。

【0039】

ファイルシステム 320 は、例えばブロック単位のディスク上フォーマット表現を有する WAF 1 ファイルシステム、例えば 4 キロバイト（kB）ブロックを使用し、`inode` を用いてファイルを表現するような WAF 1 ファイルシステムを実施する。WAF 1 ファイルシステムは、ファイルを用いてファイルシステムのレイアウトを表すメタデータを格納する。それらのメタデータ・ファイルには、特に、`inode` ファイルが含まれる。ファイルハンドル、すなわち `inode` 番号を含む識別子を用いて、ディスクから `inode` を取得する。`inode` ファイルを有するファイルシステムの構造に関する説明については、1998 年 10 月 6 日に David Hitz 他に付与された「Method for Maintaining Consistent States of a File System and for Creating User Accessible Read-Only Copies of a File System」と題する米国特許第 5,819,292 号に記載されており、この特許は、ここで参照することにより、完全に説明されたものとして本明細書に取り込まれる。

【0040】

大まかに言えば、ファイルシステムの全ての `inode` は、`inode` ファイルとして編成される。ファイルシステム（FS）情報ブロックは、ファイルシステム内の情報のレイアウトを指定するためのものであり、ファイルシステム内の他の全ての `inode` を含むファイルの `inode` を有する。各ボリュームは FS 情報ブロックを有し、FS 情報ブロックは、例えばファイルシステムの RAID グループ内の固定位置に格納することが望ましい。ルート FS 情報ブロックの `inode` は、`inode` ファイルの直接参照（指し示す）ブロックであってもよいし、`inode` ファイルの間接参照ブロック、すなわち、

次いで `inode` ファイルを直接参照するブロックであってもよい。`inode` ファイルの各直接参照ブロックは埋め込み `inode` であり、それらの各々が、間接ブロックをさらに参照し、次いでファイルまたは `vdisk` のデータブロックを参照する場合がある。

【0041】

本発明の一態様によれば、ファイルシステムは、`vdisk322` 並びにファイル `324` およびディレクトリ (`dir326`) に対するアクセス処理と同時に、ボリューム `150` および / または `qtree328` のようなストレージ・ユニットの汎用記憶空間管理を実施する。`qtree328` は、物理ボリュームの名前空間内の論理サブボリュームの属性を有する特別なディレクトリである。各ファイルシステム・ストレージ・オブジェクト (ファイル、ディレクトリ、または `vdisk`) は、例えば1つの `qtree`、クォータ、セキュリティ・プロパティに関連付けられ、他の項目は、`qtree` ごとに割り当てられる。`vdisk` およびファイル/ディレクトリは、`qtree328` の一番上の層に配置され、次いでそれらは、ファイルシステム「仮想化」層 `320` による抜粋に従って、ボリューム `150` の一番上の層に配置される。

10

【0042】

ファイルシステム `320` 内の `vdisk` ストレージ・オブジェクトは、マルチプロトコル・ストレージ・アプライアンスの `SAN` デプロイメントに関連する一方、ファイル・ストレージ・オブジェクトおよびディレクトリ・ストレージ・オブジェクトは、アプライアンスの `NAS` デプロイメントに関連する点に注意して欲しい。それらのファイルおよびディレクトリは一般に、`FC` アクセス・プロトコルや `SCSI` ブロック・アクセス・プロトコルによってアクセスすることができない。しかしながら、ファイルを `vdisk` にコンバートしてから、`SAN` プロトコルや `NAS` プロトコルによってアクセスすることは可能である。`vdisk` は、`SAN` (`FC` および `SCSI`) プロトコルによって `LUN` としてアクセスされる場合もあれば、`NAS` (`NFS` および `CIFS`) プロトコルによってファイルとしてアクセスされる場合もある。

20

【0043】

本発明の他の実施形態において、仮想化システム `300` は、ファイルシステム `320` による汎用空間管理に関し、`SAN` ストレージ・オブジェクトと `NAS` ストレージ・オブジェクトの共存が可能な仮想記憶空間を提供する。その目的のために、仮想化システム `300` は、ファイルシステムが本来備えている能力を含む、ファイルシステムの特性を利用して、ディスクを集め、それらを1つのストレージ・プールに抽出する。例えば、システム `300` は、ファイルシステムのボリューム管理機能を利用して、ディスク `130` の集まりを汎用記憶空間のプールを表す1以上のボリューム `150` として構成する。次いで、`vdisk322` およびファイル `324` をそれぞれ作成することにより、その汎用記憶プールを `SAN` デプロイメントと `NAS` デプロイメントの両方にとって利用可能にする。`SAN` および / または `NAS` デプロイメントを展開したときに、`vdisk` およびファイルは、同一の汎用記憶空間を共有するだけでなく、同一の空き記憶プールを利用する。従来のシステムとは違い、マルチプロトコル・ストレージ・アプライアンスの汎用記憶空間内には、ディスクの物理的区分が無い。

30

【0044】

マルチプロトコル・ストレージ・アプライアンスは、ユーザが単一のストレージ・リソース・プールを用いて `NAS` ストレージ・オブジェクトと `SAN` ストレージ・オブジェクトの両方を管理できるようにすることで、汎用記憶空間の管理を実質的に単純化する。特に、`SAN` デプロイメントと `NAS` デプロイメントの両方について、汎用フリー・プールからの空きブロック空間が、細粒ブロックごとに管理される。仮にそれらのストレージ・オブジェクトを個別 (別々) に管理したとすれば、ユーザは、例えば業務目的などに応じて、目的の各タイプについて特定量の「予備」ディスクを保持しなければならないであろう。そのような個別アプローチを維持するために必要とされるオーバーヘッドは、業務命令に応じて拡張可能な単一グループの予備ディスクだけを用いて、それらのオブジェクトを単一のリソースプールから管理できる場合に比べて、大きくなるであろう。`vdisk` 処

40

50

理によって個別に開放されたブロックは、N A Sオブジェクトによって即座に再使用することができる（そして、その逆も可）。そのような管理の詳細を管理者が意識する必要はない。これは、統合型マルチプロトコル・ストレージ・アプライアンスの「所有権の総コスト」利点を表す。

【0045】

仮想化システム300は更に、マルチプロトコル・ストレージ・アプライアンスの汎用記憶空間内に共存するそれらのS A Nストレージ・オブジェクトおよびN A Sストレージ・オブジェクトに対し、信頼性保証を提供する。特に、従来のS A Nシステムにおいて物理ブロックレベルで実施されているR A I Dやミラーリングのような技術による、ディスク故障に直面したときの信頼性保証は、アプライアンス100のファイルシステム320から継承された特徴である。これによって、管理者は、ファイルシステム内のv d i s kとN A Sオブジェクトに等しく使用される基礎冗長物理ストレージの全体的判断を行なうことが可能となるため、管理が簡単になる。

10

【0046】

上記のように、ファイルシステム320は、情報をファイル、ディレクトリ、およびv d i s kとしてd i s k 130のボリューム150内に編成する。各ボリューム150の基礎となるのは、ボリューム内のディスク故障（複数の場合もあり）に対する保護および信頼性を提供するR A I Dグループ140～144の集まりである。マルチプロトコル・ストレージ・アプライアンスによってサービス提供される情報は、例えばR A I D 4構成に従って保護される。このレベルの保護は、例えばアプライアンス・プラットフォーム上の同期ミラーリングにまで拡張される場合もある。ボリューム150に対し後者の保護が指定された場合、ボリューム上に作成され、R A I D 4によって保護されるv d i s k 322は、さらに同期ミラーリング保護を「継承」する。この場合、同期ミラーリング保護は、v d i s kの特徴ではなく、基礎ボリュームの特徴およびファイルシステム320の信頼性保証である。マルチプロトコル・ストレージ・アプライアンスのこの「継承」特性によれば、システム管理者が信頼性問題を扱う必要がなくなるので、v d i s kの管理は簡単になる。

20

【0047】

さらに、仮想化システム300は、ディスク130の物理構成に関するユーザの知識を必要とせず、ディスク130の帯域幅を集める。ファイルシステム320は、データが格納されるボリュームの全ディスクの帯域を集める入出力（I/O）記憶処理に従って、データを長い連続的なストライプとしてそれら複数のディスクにわたってディスクに書き込む（格納する）。ディスクに対して情報を格納または取り出しする場合、ディスクに対するそれらのI/O処理は、ユーザによって指定されるものではない。そうではなく、それらの処理はユーザにとって透過的である。なぜなら、ファイルシステムが、w r i t e a n y w h e r eレイアウト・ポリシーに従って、そのデータをボリュームの全ディスクにわたって信頼性の高い態様に「ストライプ化」するからである。ブロック・ストレージの仮想化によって、v d i s kに対するI/O帯域幅をv d i s kのサイズとは無関係にファイルシステムの基礎物理ディスクの最大帯域幅にすることが可能になる（従来のブロックアクセス製品における一般的なL U Nの物理実施形態とは異なる）。

30

40

【0048】

さらに、仮想化システムは、ファイルシステム配置ポリシー、管理ポリシーおよびブロック管理ポリシーに影響を与え、マルチプロトコル・ストレージ・アプライアンス内でv d i s kを正しく機能させる。v d i s kブロック配置ポリシーは基礎となる仮想化ファイルシステムの機能であり、変更直面したときに、ファイルシステム・ブロックとS C S I論理ブロックアドレスとの間に、永久的な結び付けは何もない。v d i s kは、透過的に再構成され、おそらくデータ・アクセス・パターン動作を変更する。

【0049】

S A NデプロイメントとN A Sデプロイメントのどちらの場合でも、ブロック割当てポリシーは、ディスクの物理特性（例えば、幾何配置、サイズ、シリンダ、セクタサイズ）

50

とは無関係である。本発明によれば、ファイルシステムは、ボリューム150内に存在するファイル324、ディレクトリ326およびvdisk322に関するファイル単位の管理を提供する。マルチプロトコル・ストレージ・アプライアンスに取り付けられたアレイにディスクを追加する場合、そのディスクは既存のボリュームに組み込まれ、ボリューム空間全体を増加させる。その空間は、任意の目的に、例えば更なるvdiskまたは更なるファイルのために使用される。

【0050】

統合型マルチプロトコル・ストレージ・アプライアンス100の管理は、システム管理者にとって利用可能なUI350およびvdiskコマンド群を使用することにより、簡単になる。UI350は例えば、コマンド・ライン・インタフェース(CLI)とグラフィカル・ユーザ・インタフェース(GUI)との両方を含む。UIは、vdiskコマンド群の実行に使用され、とりわけ、vdiskの作成、vdiskのサイズの増減、および/または、vdiskの破壊に使用される。アプライアンス100の仮想化記憶空間の特徴によれば、破壊されたvdiskの記憶空間は、例えばNAS系ファイルのために再利用することができる。vdiskは、そのアプリケーション・データに対するブロックおよびNASマルチプロトコル・アクセスを維持したまま、ユーザ制御によって増加(「拡大」)または減少(「縮小」)させることができる。

【0051】

UI350は、vdiskを作成する際に使用すべきディスクをシステム管理者が明示的に構成および指定する必要性を無くすことにより、マルチプロトコルSAN/NASストレージ・アプライアンスの管理を簡単にする。例えばvdiskを作成する場合、システム管理者は、例えばCLI352またはGUI354を通じて、単にvdisk(「LUN作成」)コマンドを発行するだけでよい。vdiskコマンドは、vdiskの作成を指示するとともに、そのvdiskの所望のサイズおよびそのvdiskのパス記述子(パス名)を指定する。これに回答して、ファイルシステム320はvdiskモジュール330と協働し、基礎ディスクによって提供された記憶空間を「仮想化」し、作成コマンドの指示通りにvdiskを作成する。具体的には、vdiskモジュール330は、vdiskコマンドを処理し、ファイルシステム320内の原始的な処理(「プリミティブ」を「呼び出し」、vdisk(LUN)の高レベル概念を実施する。例えば、「LUN作成コマンド」は、vdiskを正しい情報およびサイズで、正しい位置に作成する一連のファイルシステム・プリミティブに変換される。それらのファイルシステム・プリミティブには、ファイルinodeを作成する処理(create_file)、ストリームinodeを作成する処理(create_stream)、情報をストリームinodeに格納する処理(stream_write)などが含まれる。

【0052】

LUN作成コマンドの結果、指定サイズのvdiskが作成され、明示的に指定しなくとも、そのvdiskはRAIDによって保護される。マルチプロトコル・ストレージ・アプライアンスのディスク上のストレージ情報はタイプされない。「未加工」のビットだけがディスク上に格納される。ファイルシステムは、それらのビットをボリューム内の全ディスクにわたって、vdiskおよびRAIDグループとして編成する。従って、作成されるvdisk322を明示的に構成する必要はない。なぜなら、仮想化システム300は、vdiskをユーザにとって透過的な態様で作成するからである。作成されたvdiskは、ファイルシステムによって作成された基礎ボリュームの信頼性や記憶帯域幅などの高性能の特徴を継承する。

【0053】

また、CLI352および/またはGUI354は、vdiskモジュール330と対話し、属性や、作成されたvdiskに番号を割り当てる永久LUNマップ・バインディングを導入する。その後、それらのLUNマップ・バインディングは、ディスクのエクスポートに使用され、クライアントに対して特定のSCSI識別子(ID)として使用される。特に、作成されたvdiskは、LUNマッピング技術によってエクスポートするこ

10

20

30

40

50

とができ、SANクライアントはディスクを「見る」(アクセスする)ことができる。SAN環境では一般にvdisk(LUN)に対するアクセスを厳密に制御しなければならないため、SAN環境におけるLUNの共有は通常、クラスタ化ファイルシステム、クラスタ化オペレーティング・システムおよびマルチパス・構成などの限られた環境でしか行なわれない。マルチプロトコル・ストレージ・アプライアンスのシステム管理者は、どのvdisk(LUN)がSANクライアントにエクスポート可能であるかを判断する。vdiskをLUNとしてエクスポートした後、クライアントは、SANネットワークを介してFCPやiSCSIなどのブロックアクセスプロトコルを用いて、そのvdiskにアクセスすることができる。

【0054】

SANクライアントは通常、論理番号すなわちLUNを用いてディスクを識別および宛先指定する。しかしながら、マルチプロトコル・ストレージ・アプライアンスの「管理の容易性」特徴は、システム管理者がvdiskを論理名で管理し、宛先指定できるようにする。その目的のために、マルチプロトコル・ストレージ・アプライアンスのvdiskモジュール330は、論理名をvdiskにマッピングする。例えばvdiskを作成する場合、システム管理者はそのvdiskに「正しいサイズ」を割り当て、そのvdiskに、そのvdiskの目的とする用途を一般的に表すような名前(例えば、データベースを保持する場合、/vol/vol0/databaseなど)を割り当てる。管理インタフェースは、ストレージ・アプライアンスからクライアントにエクスポートされるLUN/vdisk(およびファイル)を名前によって管理し、それによって、ブロック単位(およびファイルベース)のストレージに対する一貫性のある統一された名前付け方式を提供する。

【0055】

マルチプロトコル・ストレージ・アプライアンスは、イニシエータ・グループ(iGROUP)を用いて論理名によるvdiskのエクスポート制御を管理する。iGROUPは、1以上のイニシエータ(クラスタ化環境が構成されているか否かに応じて異なる)に関連する1以上の宛先に割り当てられた論理名付きのエンティティである。「iGROUP作成」コマンドは、それらの宛先を実質的に「バインド」する(関連付ける)コマンドである。宛先には、論理名またはiGROUPに対するWWNアドレスやiSCSI IDが含まれる。そして、「LUNマップ」コマンドを用いて、1以上のvdiskをigroupにエクスポートする。すなわち、vdiskがigroupから「見える」ようにする。その意味では、「LUNマップ」コマンドは、NFSのエクスポートや、CIFSのシェアと同等のコマンドである。従って、WWNアドレスまたはiSCSI IDは、LUNマップコマンドによって指定された、それらのディスクに対するアクセスが許可されているクライアントを識別する。その後、ストレージ・オペレーティング・システムの内部の全ての処理に、論理名が使用される。この論理名の概念は、ユーザとマルチプロトコル・ストレージ・アプライアンスとの間の対話を含む、vdiskコマンド群全体に行き渡る。特に、以後のエクスポート処理にはすべて、igroup名変換が使用され、種々のSANクライアントにエクスポートされたLUNのリストが作成される。

【0056】

図4は、マルチプロトコル・ストレージ・アプライアンスに格納された情報をSANネットワークを介してアクセスするときに必要とされる各ステップのシーケンスを示す略フロー図である。ここで、クライアントは、ストレージ・アプライアンス100に接続されたネットワークを介し、ブロック・アクセス・プロトコルを用いてアストレージ・アプライアンス100と通信するものとする。クライアント160が、Windows(登録商標)・オペレーティング・システムを実行しているクライアント160aである場合、ネットワーク185を介して例えばブロック・アクセス・プロトコルが使用される。一方、クライアントがUNIX(登録商標)・オペレーティング・システムを実行しているクライアント160bである場合、ネットワーク165を介して例えばiSCSIプロトコルが使用される。シーケンスはステップ400で始まり、ステップ402へ進み、そこで、クライアントは、マルチプロトコル・ストレージ・アプライアンス内の情報にアクセスする

10

20

30

40

50

ための要求を生成、ステップ404で、その要求が、ネットワーク185、165を介して、従来のFCP・プロトコルまたはiSCSI・プロトコルのアクセス要求として転送される。

【0057】

ステップ406で、その要求はストレージ・アプライアンス100のネットワーク/アダプタ126、125によって受信され、そこで統合ネットワーク・プロトコル・スタックによって処理され、ステップ480で、それが仮想化システム300に渡される。具体的には、要求がFCP要求であった場合、その要求は、例えばFCドライバ230によってデータをアクセス(すなわち、読み書き)する4kブロック要求として処理される。要求がiSCSIプロトコル要求であった場合、その要求は、メディア・アクセス層(Intel ギガビット・イーサネット(登録商標))で受信され、TCP/IPネットワーク・プロトコル層を通して仮想化システムに渡される。

10

【0058】

SCSIプロトコルに関連するアドレス情報を含むコマンドおよび制御処理は一般に、ディスクまたはLUNに対して行なわれる。しかしながら、ファイルシステムはLUNを認識しない。従って、その要求に含まれるSCSIコマンドにตอบสนองするために、仮想化システムのSCSIターゲット・モジュール310は、LUNのエミュレーションを開始する(ステップ410)。その目的のために、SCSIターゲット・モジュールは、SCSIプロトコルを利用する一連のアプリケーション・プログラム・インタフェース(API 360)であって、iSCSIドライバ228とFCドライバ230の両方に対する一貫性のあるインタフェースを実施する一連のアプリケーション・プログラム・インタフェースを有する。SCSIターゲット・モジュールは更に、LUNをvdiskに実質的に変換するマッピング/変換手順を実施する。ステップ412で、SCSIターゲット・モジュールは、その要求に関する例えばFCルーティング情報などのアドレス情報をファイルシステムの内部構造にマッピングする。

20

【0059】

ファイルシステム320は例えばメッセージを利用したシステムである。従って、SCSIターゲット・モジュール310は、SCSI要求をファイルシステムに対する処理を表すメッセージに変換する。SCSIターゲット・モジュールによって生成されるメッセージには、例えば、ファイルシステム上に表現されたvdiskオブジェクトのパス名(例えば、パス記述子など)およびファイル名(例えば、特別なファイル名など)を伴う処理のタイプ(例えば、読出し、書込みなど)がある。SCSIターゲット・モジュールは、メッセージを例えばファンクションコール365としてファイルシステム320に渡し、そこでその処理が実施される。

30

【0060】

ファイルシステム320は、メッセージの受信にตอบสนองして、パス名をinode構造にマッピングし、vdisk322に対応するファイル・ハンドルを得る。ファイル・ハンドルを得た後、ストレージ・オペレーティング・システム200は、そのハンドルをディスク・ブロックに変換する。すなわち、そのブロック(inode)をディスクから取り出す。大まかに言えば、ファイル・ハンドルとは、ファイルシステムの内部で使用されるデータ構造の内部表現、すなわちinodeデータ構造の内部表現である。ファイル・ハンドルは通常、ファイルID(inode番号)、スナップショットID、生成ID、およびフラグを含む複数の構成要素からなる。ファイルシステムは、ファイル・ハンドルを用いて、特別なファイルinodeと、ディスク130上で実施されるファイルシステム構造内のvdiskを含む少なくとも1つの関連ストリーム・inodeとを取得する。

40

【0061】

要求されたデータがコア内に、すなわちメモリ134内に無かった場合、ステップ414で、ファイルシステムは、要求されたデータをディスク130から読み込む(取り出す)処理を生成する。情報がメモリ内に無かった場合、ファイルシステム320は、inodeファイル内の索引としてinode番号を使用し、適当なエントリにアクセスし、論

50

理ボリューム・ブロック番号 (VBN) を取り出す。次にファイルシステムは、その論理 VBN をディスク・ストレージ (RAID) 層 240 に渡し、そこで、その論理番号をブロック番号にマッピングし、それをディスク・ドライバ層 250 の適当なドライバ (例えば、SCSI) に送る。ディスク・ドライバは、ディスク 130 からそのディスク・ブロック番号にアクセスし、それをメモリ 124 内の要求データブロック (複数の場合もあり) に読み込む。ステップ 416 で、要求されたデータは仮想化システム 300 によって処理される。例えば、そのデータは、vdisk に対する読出しまたは書込み処理に関連して、または、vdisk に対する照会コマンドに関連して処理される。

【0062】

SCSI ターゲット・モジュールは、要求された vdisk に関する重要な意味を「シミュレート」することによって、従来の SCSI プロトコルのサポートをエミュレートする。そのような情報は、SCSI ターゲット・モジュールによって計算することもできるし、例えば vdisk の属性ストリーム inode に永久的に格納しておいてもよい。ステップ 418 で、SCSI ターゲット・モジュール 310 は、要求されたブロック単位の情報 (ファイルシステム 320 によって提供されたファイル単位の情報から変換されるようなもの) をブロック・アクセス (SCSI) ・プロトコル・メッセージの中にロードする。例えば、SCSI ターゲット・モジュール 310 は、SCSI 照会コマンド要求に回答して、vdisk のサイズ等の情報を SCSI プロトコル・メッセージの中にロードすることができる。要求が完了すると、ストレージ・アプライアンス (およびオペレーティング・システム) は、ネットワークを介して回答 (例えば、SCSI 「容量」 応答メッセージ) をクライアントに返す (ステップ 420)。そして、シーケンスはステップ 422 で終了する。

【0063】

マルチプロトコル・ストレージ・アプライアンスで受信されたクライアント要求について、データ記憶アクセスを実施するために必要とされる、上で説明したストレージ・オペレーティング・システム層を通るソフトウェア「パス」は、代替としてハードウェアで実施してもよいことに注意して欲しい。すなわち、本発明の代替実施形態において、オペレーティング・システム層 (仮想化システム 300 を含む) を通るストレージ・アクセス・要求データ・パスは、フィールド・プログラマブル・ゲート・アレイ (FPGA) または特定用途向け IC (ASIC) の中に論理回路として実施することもできる。この種のハードウェア実施形態は、クライアント 160 によって発行されたファイル・アクセス要求やブロック・アクセス要求に回答してアプライアンス 100 により提供されるストレージ・サービスの能力を向上させる。さらに、本発明の他の実施形態において、ネットワークやストレージ・アダプタ 125 ~ 128 などの処理要素は、プロセッサ 122 からパケット処理やストレージ・アクセス処理などの負荷の全部または一部を分担することにより、マルチプロトコル・ストレージ・アプライアンスによって提供されるストレージ・サービスの能力を向上させるように構成される場合がある。当然ながら、本明細書に記載される種々の処理、アーキテクチャおよび手順は、ハードウェアで実施することも、ファームウェアで実施することも、あるいはソフトウェアで実施することも可能である。

【0064】

有利なことに、統合型マルチプロトコル・ストレージ・アプライアンスは、データの完全性を確保しつつ、アクセス制御を行なうことができ、必要であれば、全てのプロトコルについてファイルおよび vdisk の共有を行なうこともできる。さらに、ストレージ・アプライアンスは、埋め込み型 / 統合型の仮想化機能を備え、NAS ストレージ・オブジェクトおよび SAN ストレージ・オブジェクトを作成するとき、ユーザがストレージ・リソースを分配する必要性が無い。それらの機能は、アプライアンス内の汎用空間管理に関し、SAN ストレージ・オブジェクトと NAS ストレージ・オブジェクトの共存を可能にする仮想化記憶空間を含む。さらに、統合型ストレージ・アプライアンスは、同一の vdisk に対するブロック・アクセス・プロトコル (iSCSI や FCP) の同時サポートを提供するだけでなく、クラスタリングをサポートする異種 SAN 環境も提供する。要

10

20

30

40

50

するに、マルチプロトコル・ストレージ・アプライアンスは、全てのストレージ・アクセス・プロトコルについて、単一の統一されたストレージ・プラットフォームを提供する。

【0065】

上記の説明は、本発明の特定の実施形態に関するものである。しかしながら、記載した実施形態に対して他の変更および修正を施し、それらの実施形態の利点のうちの一部または全部が得られるようにすることも可能であることは、明らかである。例えば、本発明の教示は、コンピュータ上で実行されるプログラム命令を有するコンピュータ読取可能媒体も含めて、ソフトウェアで実施することも、ハードウェアで実施することも、ファームウェアで実施することもでき、それらの組み合わせで実施することもできるものと当然ながら考えられる。従って、この説明は、単なる例として解釈すべきものであり、発明の範囲を制限するものとして解釈してはならない。従って、特許請求の範囲のねらいは、それらの変更および修正が、本発明の真の思想および範囲に含まれるようにカバーすることにある。

10

【図面の簡単な説明】

【0066】

【図1】ストレージ・エリア・ネットワーク（SAN）環境およびネットワーク・アタッチド・ストレージ（NAS）環境で動作するように構成された本発明によるマルチプロトコル・ストレージ・アプライアンスを示す略ブロック図である。

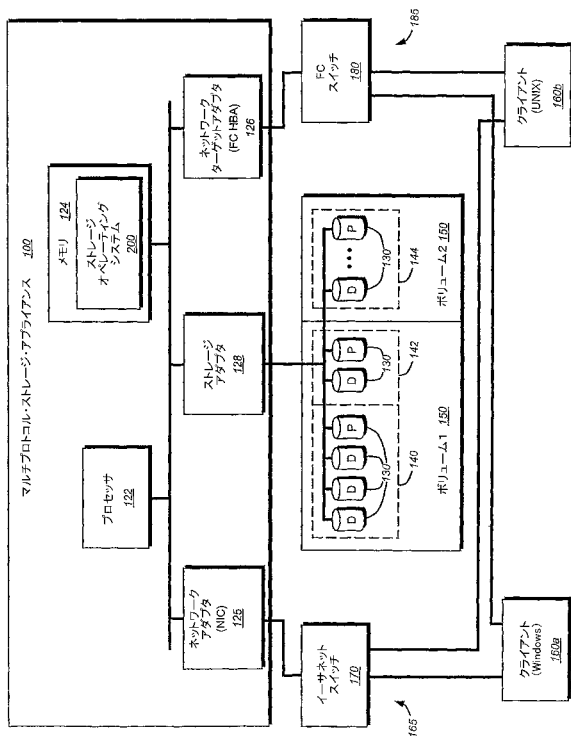
【図2】本発明で使用するのに都合がよいマルチプロトコル・ストレージ・アプライアンスのストレージ・オペレーティング・システムを示す略ブロック図である。

20

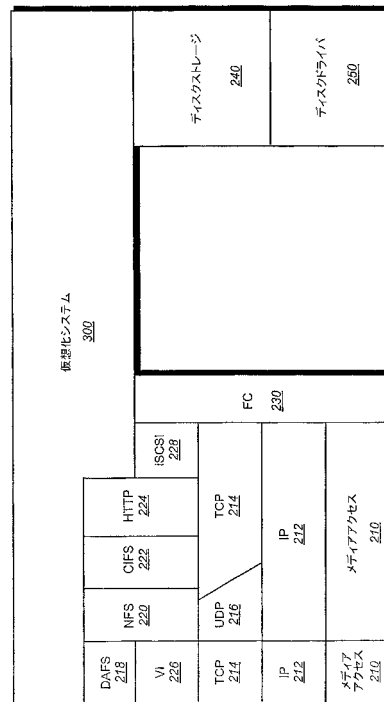
【図3】ファイルシステムによって実施され、仮想化モジュールに作用する本発明による仮想化システムを示す略ブロック図である。

【図4】SANネットワークを介してマルチプロトコル・ストレージ・アプライアンスに格納された情報にアクセスするときに必要なとされる各ステップのシーケンスを示すフロー図である。

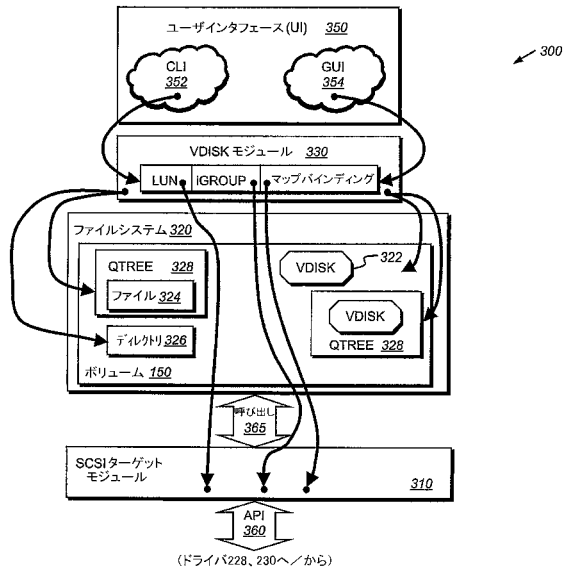
【図1】



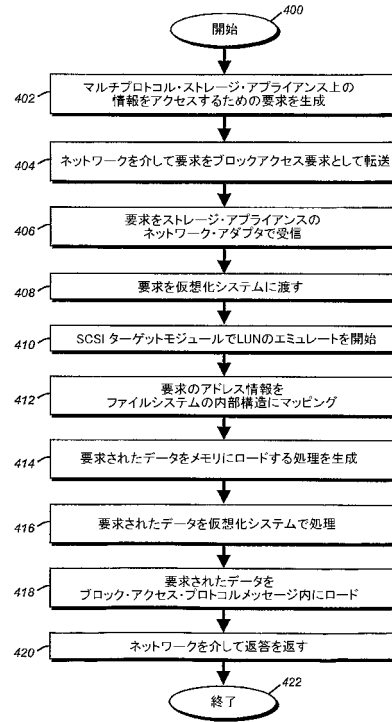
【図2】



【 図 3 】



【 図 4 】



【 手続 補正 書 】

【 提出 日 】 平成 17 年 6 月 22 日 (2005.6.22)

【 手続 補正 1 】

【 補正 対 象 書 類 名 】 特 許 請 求 の 範 囲

【 補正 対 象 項 目 名 】 全 文

【 補正 方 法 】 変 更

【 補正 の 内 容 】

【 特 許 請 求 の 範 囲 】

【 請 求 項 1 】

ネットワーク・アタッチド・ストレージ (NAS) デプロイメントと、ストレージ・エリア・ネットワーク (SAN) デプロイメントとの両方に対し、記憶装置に格納された情報に対するファイル・プロトコル・アクセスおよびブロック・プロトコル・アクセスを提供するように構成されたマルチプロトコル・ストレージ・アプライアンスであって、

仮想化モジュールと協働し、前記記憶装置によって提供される記憶空間を仮想化するように構成されたストレージ・オペレーティングシステムであって、ファイル・プロトコル・アクセス要求とブロック・プロトコル・アクセス要求の両方からの情報を前記記憶空間に統一した態様で格納する、ストレージ・オペレーティング・システムを含む、マルチプロトコル・ストレージ・アプライアンス。

【 請 求 項 2 】

前記ファイルシステムは、情報をファイル、ディレクトリおよび仮想ディスク (vdisk) として論理編成し、ファイルおよびディレクトリに対するファイル単位のアクセスを可能にするとともに、vdiskに対するブロック単位のアクセスをさらに可能することにより、NASアプライアンスとSANアプライアンスの統一した記憶方法を提供し、前記vdiskのそれぞれは、エミュレート・ディスクとして解釈される特別なファイルタイプである、請求項1に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 3】

前記仮想化モジュールは、v d i s kモジュールと、S C S I (Small Computer Systems Interface) ターゲット・モジュールとを含む、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 4】

前記 v d i s kモジュールは、前記ファイルシステム上に層として形成され、システム管理者がマルチプロトコル・ストレージ・アプライアンスに対して発行したコマンドに回答して、管理インタフェースによるアクセスが有効にされる、請求項 3 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 5】

前記管理インタフェースは、ユーザ・インタフェース (U I) を含む、請求項 4 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 6】

前記 v d i s kモジュールは、前記 U I を通じて発行された一連の v d i s k コマンドを実行することにより、S A N デプロイメントを管理する、請求項 5 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 7】

前記 v d i s k コマンドは、前記ファイルシステムおよび前記 S C S I ターゲット・モジュールに作用し、前記 v d i s k を実施する原始的なファイルシステム処理に変換される、請求項 6 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 8】

前記 S C S I ターゲット・モジュールは、論理ユニット番号 (L U N) を v d i s k に変換するマッピング手順を提供することにより、ディスクまたは前記 L U N のエミュレーションを開始する、請求項 7 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 9】

前記 S C S I ターゲット・モジュールは、S A N ブロック空間とファイルシステム空間との間の変換層を提供する、請求項 8 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 10】

前記ファイルシステムによる汎用記憶空間管理に関し、前記仮想化記憶空間は、S A N ストレージ・オブジェクトと N A S ストレージ・オブジェクトを共存させることが可能である、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 11】

前記ファイルシステムは前記仮想化モジュールと協働し、前記仮想化記憶空間内に共存する前記 S A N ストレージ・オブジェクトおよび前記 N A S ストレージ・オブジェクトに信頼性保証を与える、請求項 10 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 12】

前記ファイルシステムは、前記記憶装置に格納された情報に対するブロック単位のアクセスをする際に使用されるボリューム管理機能を提供する、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 13】

前記記憶装置はディスクである、請求項 12 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 14】

前記ファイルシステムは、(i) ストレージ・オブジェクトの名前付けなどのファイルシステム・セマンティック、および、(ii) ボリューム・マネージャに関連する機能を提供する、請求項 1 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 15】

前記ボリューム・マネージャに関連する機能は、

前記記憶装置を集める機能と、

前記記憶装置の記憶帯域幅を集める機能と、

ミラーリングまたは R A I D (Redundant Array of Independent Disks) のような信頼性保証機能と

のうちの少なくとも 1 つを含む、請求項 1 4 に記載のマルチプロトコル・ストレージ・アプライアンス。

【請求項 1 6】

マルチプロトコル・ストレージ・アプライアンスに接続された記憶装置上での情報の編成に関するストレージ・サービスを提供するように構成されたマルチプロトコル・ストレージ・アプライアンスのストレージ・オペレーティング・システムであって、

ブロック単位のアクセス・プロトコルまたはファイル単位のアクセス・プロトコルを用いて前記マルチプロトコル・ストレージ・アプライアンスに格納された情報にアクセスするクライアントに対し、データ・パスを提供する統合型ネットワーク・プロトコル・スタックと、

仮想化モジュールと協働し、前記記憶装置によって提供される記憶空間を仮想化するファイルシステムと、

からなるストレージ・オペレーティング・システム。

【請求項 1 7】

前記ファイルシステムは、情報をファイル、ディレクトリおよび仮想ディスク (v d i s k) として論理編成することにより、ファイル単位のアクセス・プロトコルとブロック単位のアクセス・プロトコルを用いた、ネットワーク・アタッチド・ストレージ (N A S) アプライアンスとストレージ・エリア・ネットワーク (S A N) アプライアンスの統一的記憶方法を提供し、前記 v d i s k のそれぞれが、エミュレート・ディスクとして解釈される特別なファイルタイプである、請求項 1 6 に記載のストレージ・オペレーティング・システム。

【請求項 1 8】

前記ブロック単位のアクセス・プロトコルは、「 S C S I e n c a p s u l a t e d o v e r T C P (i S C S I) 」および「 S C S I e n c a p s u l a t e d o v e r F C (F C P) 」のような S C S I (Small Computer Systems Interface) 系のプロトコルを含む、請求項 1 7 に記載のストレージ・オペレーティング・システム。

【請求項 1 9】

前記統合型ネットワーク・プロトコル・スタックは、

ネットワーク・プロトコル層と、

前記ネットワーク・プロトコル層に接し、前記ファイルシステムにより編成されたファイルおよびディレクトリに対するファイル単位のプロトコル・アクセスを提供するファイルシステム・プロトコル層と、

前記ネットワーク・プロトコル層の上に配置され、前記ファイルシステムによって編成された v d i s k に対するブロック単位のプロトコル・アクセスを提供する i S C S I ドライバと、

を含む、請求項 1 8 に記載のストレージ・オペレーティング・システム。

【請求項 2 0】

前記統合型ネットワーク・プロトコル・スタックは、前記ファイルシステム・プロトコル層のファイル・アクセス・プロトコルのための直接アクセス搬送機能を提供する仮想インタフェース層をさらに含む、請求項 1 9 に記載のストレージ・オペレーティング・システム。

【請求項 2 1】

前記統合型ネットワーク・プロトコル・スタックは、前記ファイルシステムによって編成された v d i s k に対するブロックアクセス要求を送受信するように構成されたファイバ・チャンネル (F C) ドライバをさらに含む、請求項 2 0 に記載のストレージ・オペレーティング・システム。

【請求項 2 2】

前記 F C ドライバおよび前記 i S C S I ドライバは、F C 特有のアクセス制御および i S C S I 特有のアクセス制御を v d i s k に提供し、マルチプロトコル・ストレージ・アプライアンス上の v d i s k にアクセスするときの、i S C S I および F C P に対する v d i s k のエクスポートの管理をさらに提供する、請求項 2 1 に記載のストレージ・オペレーティング・システム。

【請求項 2 3】

ネットワーク・アタッチド・ストレージ (N A S) デプロイメントと、ストレージ・エリア・ネットワーク (S A N) デプロイメントとの両方に対し、マルチプロトコル・ストレージ・アプライアンスの記憶装置に格納された情報にアクセスするためのファイル・プロトコル・アクセスおよびブロック・プロトコル・アクセスを提供する方法であって、格納された情報を保持するための単一の記憶空間を提供するステップと、

(i) 前記アプライアンスを第 1 のネットワークに接続するネットワーク・アダプタと、(i i) N A S クライアントが前記格納された情報をファイルとしてアクセスするために発行するファイル単位の要求に対し、前記アプライアンスが応答できるようにするためのファイルシステム機能とを用いて、N A S サービスを提供するステップと、

(i) 前記アプライアンスを第 2 のネットワークに接続するネットワーク・ターゲット・アダプタと、(i i) S A N クライアントが前記格納された情報を仮想ディスク (v d i s k) としてアクセスするために発行するブロック単位の要求に対し、前記アプライアンスが応答できるようにするためのボリューム管理機能とを用いて、S A N サービスを提供するステップと、

からなる方法。

【請求項 2 4】

前記マルチプロトコル・ストレージ・アプライアンスに格納されたファイルおよび v d i s k の名前による管理を提供することにより、ファイル単位のストレージおよびブロック単位のストレージのための一様な名前付け方式を提供するステップと、

前記記憶装置に格納された名前付きファイルおよび v d i s k の階層構造を提供するステップと、

をさらに含み、前記 v d i s k のそれぞれが、エミュレート・ディスクとして解釈される特別なファイルタイプである、請求項 2 3 に記載の方法。

【請求項 2 5】

ネットワーク・アタッチド・ストレージ (N A S) デプロイメントと、ストレージ・エリア・ネットワーク (S A N) デプロイメントとの両方に対し、マルチプロトコル・ストレージ・アプライアンスの記憶装置に格納されたストレージ・オブジェクトにアクセスするためのファイル・プロトコル・アクセスおよびブロック・プロトコル・アクセスを提供する方法であって、

前記記憶装置を汎用記憶空間を表す 1 以上のいボリュームとして編成するステップと、

前記汎用記憶空間内での S A N ストレージ・オブジェクトと N A S ストレージ・オブジェクトの共存を許可するステップと、

前記 S A N ストレージ・オブジェクトおよび前記 N A S ストレージ・オブジェクトにアクセスするためのブロック単位の要求およびファイル単位の要求を前記ストレージ・アプライアンスのマルチプロトコル・エンジンで受信するステップと、

前記ブロック単位の要求および前記ファイル単位の要求に応答し、前記 S A N ストレージ・オブジェクトおよび前記 N A S ストレージ・オブジェクトにアクセスし、それらを返すステップと、

からなる方法。

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/US03/23597
A. CLASSIFICATION OF SUBJECT MATTER		
IPC(7) : G06F 12/08, 15/16 US CL : 709/211, 245; 711/203 According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) U.S. : 709/208, 211, 217, 219, 230, 245; 711/117, 147, 148, 202, 203		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) Please See Continuation Sheet		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X, E	US 6,606,690 B2 (PADOVANO) 12 AUGUST 2003 Abstract, Column 2, Lines 15-34, Column 3, Lines 1-39, Column 5, Lines 1-54, Column 9, Lines 38-48, Column 10, Lines 1-32, Column 22, Lines 47-62	1-27
A	CALLAGHAN, B. "NFS Version 3 Protocol Specification", Request for Comments (RFC) 1813, June 1995	1-27
A	US 6,185,655 B1 (PEPING) 06 FEBRUARY 2001 Entire document	1-27
A, P	LU, Y. "Performance Study of iSCSI-Based Storage Subsystems", IEEE Communications Magazine, pp. 76-82, August 2003	1-27
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents:		
"A"	document defining the general state of the art which is not considered to be of particular relevance	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
"E"	earlier application or patent published on or after the international filing date	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"L"	document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"O"	document referring to an oral disclosure, use, exhibition or other means	"&" document member of the same patent family
"P"	document published prior to the international filing date but later than the priority date claimed	
Date of the actual completion of the international search 26 March 2004 (26.03.2004)		Date of mailing of the international search report 14 APR 2004
Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US Commissioner for Patents P.O. Box 1450 Alexandria, Virginia 22313-1450 Facsimile No. (703) 305-3230		Authorized officer Marc D. Thompson Telephone No. 703-305-3900

INTERNATIONAL SEARCH REPORT

PCT/US03/23

Continuation of B. FIELDS SEARCHED Item 3:
EAST full text USPAT, EPO, JPO -- IEEE
Terms: storage area network (SAN), network attached storage (NAS), virtual disk/storage/volume/memory/address, file/block access, file system/volume/disk mounting, ISCSI, RDMA

フロントページの続き

(81) 指定国 AP(GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VC, VN, YU, ZA, ZM, ZW

(72) 発明者 パウロウスキー, ブライアン
アメリカ合衆国カリフォルニア州 9 4 3 0 1, パロアルト, ハミルトン・アベニュー・1 1 5 6

(72) 発明者 スリニヴァサン, モハン
アメリカ合衆国カリフォルニア州 9 5 0 5 1, サンタクララ, クロニン・ドライブ・1 1 7

(72) 発明者 リー, ハーマン
アメリカ合衆国カリフォルニア州 9 4 0 4 1, マウンテンビュー, ナンバー 1 6 0, レインボー・ドライブ・6 0 0

(72) 発明者 ラジャン, ビジャヤン
アメリカ合衆国カリフォルニア州 9 4 0 8 5 - 1 6 1 2, サニーベイル, コスタ・メサ・テラス・ナンバーエイ・4 3 1

(72) 発明者 ピットマン, ジョセフ, シー
アメリカ合衆国ノースカロライナ州 2 7 5 0 2, アベックス, オークウェル・コート・1 0 0 2

Fターム(参考) 5B014 EB04

5B065 BA01 CA30 CC10 CE30 EA02 EA31

5B082 FA07 HA08

【要約の続き】

アンスの統一的記憶方法を提供する。