

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4606455号
(P4606455)

(45) 発行日 平成23年1月5日(2011.1.5)

(24) 登録日 平成22年10月15日(2010.10.15)

(51) Int.Cl. F 1
G 0 6 F 13/10 (2006.01) G 0 6 F 13/10 3 4 0 A
G 0 6 F 3/06 (2006.01) G 0 6 F 3/06 3 0 4 R

請求項の数 7 (全 29 頁)

(21) 出願番号	特願2007-328396 (P2007-328396)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成19年12月20日(2007.12.20)		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2009-151519 (P2009-151519A)	(74) 代理人	100092152 弁理士 服部 毅巖
(43) 公開日	平成21年7月9日(2009.7.9)	(72) 発明者	田村 雅寿 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21年3月19日(2009.3.19)	(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
		(72) 発明者	荻原 一隆 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

最終頁に続く

(54) 【発明の名称】 ストレージ管理装置、ストレージ管理プログラムおよびストレージシステム

(57) 【特許請求の範囲】

【請求項1】

データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムにおけるストレージ装置を管理するストレージ管理装置において、

前記ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段と、

前記ストレージ装置における読み出し処理と書き込み処理との処理性能を示すパトロール性能を測定するパトロール性能測定手段と、

前記パトロール性能測定手段によって測定された前記パトロール性能に基づいてパトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないように、前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段と、

を有することを特徴とするストレージ管理装置。

【請求項2】

前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、

前記パトロール処理実行手段は、さらに、前記パトロール処理において、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記ストレージ装置に記憶されている前記データと、前記二重化ストレージ装置に記憶されている前記データと同一内容の前

記冗長データとを読み出し、読み出した前記データと前記冗長データとを比較することによって、比較した前記データと前記冗長データとが一致することを確認することを特徴とする請求項1記載のストレージ管理装置。

【請求項3】

前記パトロール性能測定手段は、前記ストレージ装置の記憶領域が複数に区分された領域である測定単位領域ごとに、前記パトロール性能を測定し、

前記パトロール処理実行手段は、前記パトロール処理を、前記ストレージ装置の記憶領域が前記測定単位領域からさらに複数に区分された領域である実行単位領域ごとに実行し、

前記パトロール流量規制手段は、前記実行単位領域ごとに前記パトロール処理の前記パトロール流量を規制することを特徴とする請求項1記載のストレージ管理装置。

10

【請求項4】

現在の時刻を示す時刻情報を取得する時刻情報取得手段を有し、

前記パトロール性能測定手段は、前記時刻情報取得手段によって取得された前記時刻情報が示す時刻が属する時間帯ごとに、測定した前記パトロール性能を複数種類保持し、

前記パトロール流量規制手段は、前記時刻情報取得手段によって取得された前記時刻情報が示す前記現在の時刻に基づいて、前記パトロール性能測定手段が保持する前記パトロール性能から前記現在の時刻が属する時間帯に対応する前記パトロール性能を選択し、選択した当該パトロール性能に基づいて前記パトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないように前記パトロール流量を規制することを特徴とする

20

請求項1記載のストレージ管理装置。

【請求項5】

前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、

前記パトロール処理実行手段は、前記パトロール処理として、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記二重化ストレージ装置に記憶されている前記データと同一内容の前記冗長データとを読み出し、読み出した前記冗長データを前記ストレージ装置内の前記データの記憶領域に対して上書きで書き込む処理を実行することを特徴とする請求項1記載のストレージ管理装置。

30

【請求項6】

コンピュータに、データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムにおけるストレージ装置の管理処理を実行させるストレージ管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の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段、

前記ストレージ装置における読み出し処理と書き込み処理との処理性能を示すパトロール性能を測定するパトロール性能測定手段、

前記パトロール性能測定手段によって測定された前記パトロール性能に基づいてパトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないように、前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段、

40

として機能させることを特徴とするストレージ管理プログラム。

【請求項7】

データをネットワークで接続された複数のストレージ装置に分散して記憶し、前記ストレージ装置を管理するストレージ管理装置を有するストレージシステムにおいて、

前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、

50

前記ストレージ管理装置は、
前記ストレージ装置の記憶領域が正常に稼働していることを確認するパトロール処理を実行するパトロール処理実行手段と、
前記ストレージ装置における読み出し処理と書き込み処理との処理性能を示すパトロール性能を測定するパトロール性能測定手段と、
前記パトロール性能測定手段によって測定された前記パトロール性能に基づいてパトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないように、前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段と、
を有し、

10

前記パトロール処理実行手段は、さらに、前記パトロール処理において、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記ストレージ装置に記憶されている前記データと、前記二重化ストレージ装置に記憶されている前記データと同一内容の前記冗長データとを読み出し、読み出した前記データと前記冗長データとを比較することによって、比較した前記データと前記冗長データとが一致することを確認す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】

20

本発明はストレージ管理装置、ストレージ管理プログラム、およびストレージシステムに関し、特にデータを複数のストレージノードで分散管理するストレージ管理装置、ストレージ管理プログラム、およびストレージシステムに関する。

【背景技術】

【0002】

現在、コンピュータを用いたデータ処理が広く行われるようになり、データを蓄積・利用するためのストレージ技術が一層重要となっている。この点、データアクセスの高速化およびデータ保管に対する高信頼化を実現するストレージ技術として、従来はRAID (Redundant Arrays of Independent Disks) が一般的に利用されてきた。RAIDでは、データを必要に応じて分割・複製して、複数のストレージ装置に分散して配置する。これにより、複数のストレージ装置間で負荷が分散されることによる高速化、データが冗長化されることによる高信頼化が実現されている。

30

【0003】

そして近年では、さらに高速化・高信頼化を図るため、RAIDの考え方を応用した分散ストレージシステムが構築されるようになってきている。分散ストレージシステムは、複数のストレージノードと、ストレージノード間を接続するネットワークを備えている。各ストレージノードは、内部に備えた、または外部に設けられて対応付けられたストレージ装置を管理しているとともに、ネットワーク通信機能を備えている。分散ストレージシステムは、この複数のストレージノードにデータを分散して記憶することで、システム全体として一層の高速化・高信頼化を実現している。

40

【0004】

このようなストレージシステムの有するストレージ装置が正常であるか否かを診断するためには、一般に、稼働しているかどうかの診断（いわゆる生存確認）が行われる。この生存確認による診断は、処理が比較的短時間で可能であるとともに、処理の負荷が小さいため、処理によるシステムの通常の業務のアクセスに対する影響も比較的小さいというメリットがある。

【0005】

このようなストレージシステムを始めとする、コンピュータシステムの性能の監視に関して、ストレージネットワークの構成要素から収集した性能情報に基づいて、以降の情報収集の対象範囲や程度を自動調整するシステムが知られている（例えば、特許文献1参照

50

)。

【0006】

また、ネットワークノードの負荷に応じて監視トラフィックの発生間隔を自動調整するネットワーク管理装置が知られている（例えば、特許文献2参照）。

また、管理対象である計算機の監視項目ごとの稼働性能値に基づいて、採取すべき監視項目を特定する管理システムが知られている（例えば、特許文献3参照）。

【特許文献1】特開2005-157933号公報

【特許文献2】特開平09-270794号公報

【特許文献3】特開2004-206495号公報

【発明の開示】

【発明が解決しようとする課題】

【0007】

しかし、上記の稼働しているかどうかの診断を行ったのみでは、ストレージ装置自体が稼働していることは確認できるが、ストレージ装置内部の特定の領域のみに発生した障害を検出することは不可能である。

【0008】

特に、例えばストレージ装置がRAID5によって構成されている場合のように、ストレージ装置間で冗長構成がとられている場合には、一つの領域の障害が検出されれば、これを復旧することが可能である。しかし、冗長構成によって対応付けられている二つの領域について、障害が同時に発生した場合には、障害によって失われたデータを復旧することが不可能である。これにより、ストレージ装置が稼働しているかどうかの診断だけでなく、ストレージ装置の全領域にわたって、障害が発生しているかどうかについても診断することが重要である。

【0009】

しかし、ストレージ装置が有する記憶領域の全領域について診断することは、ストレージシステムに与える負荷が大きい上に、処理の終了まで長時間を要する。このため、ストレージ装置の記憶領域の診断は、通常の業務によるアクセスへの影響が無視できなくなるという問題点がある。

【0010】

また、分散ストレージシステムにおいては、多種多様なストレージ装置が存在する場合もある。このような場合には、同じ内容の診断処理を行っても、ストレージ装置の種類によって、その所要時間や通常の業務によるアクセスへの影響が異なる場合がある。このため、ストレージ装置の全領域の診断を行った場合の通常の業務に対する影響の評価が困難であるという問題点がある。

【0011】

また、分散ストレージシステムにおいては、ストレージノード間で冗長構成がとられる場合もある。データが正常かどうかを診断するためには、冗長構成をとっている各ストレージノードのデータが正しく冗長構成を保持できているかどうかの診断を行う必要がある。ここで、この分散ストレージシステムは、多種多様な構成のストレージノードが混在することによってシステムが構築されている場合もある。この場合、ストレージ装置と同様に、同じ内容の診断処理を行っても、各ストレージノードの構成によって、その所要時間や通常の業務によるアクセスへの影響が異なる場合がある。このため、ストレージ装置の全領域の診断を行った場合の通常の業務に対する影響の評価が困難であるという問題点がある。

【0012】

しかし、特許文献1～3に記載されている技術では、システムを構成するコンピュータに対するアクセス頻度などの、システムの様々な性能情報を収集することが可能であるが、それ自体では、ストレージ装置の各領域の異常を診断することはできず、また、直接的に診断の負荷を軽減することもできない。

【0013】

10

20

30

40

50

本発明はこのような点に鑑みてなされたものであり、ストレージ装置へのアクセスに対する、パトロール処理の実行に基づく負荷による影響を抑えつつ、ストレージシステムの信頼性を高めることができるストレージ管理装置、ストレージ管理プログラムおよ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】

【0014】

この、データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムにおけるストレージ装置を管理するストレージ管理装置は、前記ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段と、前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段と、を有することを要件とする。

10

【0015】

このようなストレージ管理装置によれば、パトロール処理実行手段により、ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理が実行される。また、パトロール流量規制手段により、パトロール処理実行手段によって実行されているパトロール処理の速さであるパトロール流量が規制される。

【0016】

また、コンピュータに、データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムにおけるストレージ装置の管理処理を実行させるストレージ管理プログラムは、前記コンピュータを、前記ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段、前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段、として機能させることを要件とする。

20

【0017】

このようなストレージ管理プログラムを実行するコンピュータによれば、パトロール処理実行手段により、ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理が実行される。また、パトロール流量規制手段により、パトロール処理実行手段によって実行されているパトロール処理の速さであるパトロール流量が規制される。

【0018】

また、データをネットワークで接続された複数のストレージ装置に分散して記憶し、前記ストレージ装置を管理するストレージ管理装置を有するストレージシステムは、前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、前記ストレージ管理装置は、前記ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段と、前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段と、を有し、前記パトロール処理実行手段は、さらに、前記パトロール処理において、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記ストレージ装置に記憶されている前記データと、前記二重化ストレージ装置に記憶されている前記データと同一内容の前記冗長データとを読み出し、読み出した前記データと前記冗長データとを比較することによって、比較した前記データと前記冗長データとが一致することを確認することを要件とする。

30

40

【0019】

このようなストレージシステムによれば、管理情報記憶手段により、当該ストレージ装置に記憶されているデータと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報が記憶されている。そして、パトロール処理実行手段により、ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理が実行される。また、パトロール流量規制手段により、パトロール処理実行手段によって実行されているパトロール処理の速さであるパトロール流量が規制される。さらに、パトロール処理

50

実行手段により、パトロール処理において、管理情報記憶手段に記憶されている管理情報が参照され、ストレージ装置に記憶されているデータと、二重化ストレージ装置に記憶されているデータと同一内容の冗長データとが読み出され、読み出されたデータと冗長データとを比較することによって、比較されたデータと冗長データとが一致することが確認される。

【発明の効果】

【0020】

開示のストレージ管理装置、ストレージ管理プログラムおよびストレージシステムによれば、パトロール処理の実行に基づく負荷による、ストレージ装置へのアクセスに対する影響を抑えつつ、ストレージシステムの信頼性を高めることができる。

10

【発明を実施するための最良の形態】

【0021】

以下、本発明の実施の形態について、図面を参照して説明する。

本実施の形態のストレージ管理装置1は、多種多様なストレージ装置を管理する多種多様なストレージノードがネットワークでつながれているマルチノードストレージシステムにおいて、ストレージ装置の各記憶領域へのデータの読み書きが正常に行えるかを診断するパトロール処理を、自律的に調整して実行することで、アクセスへの影響を抑えつつ、ストレージシステムの信頼性を向上させるものである。

【0022】

図1は、本実施の形態の概要を示す図である。図1に示すストレージシステムは、データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムである。このストレージシステムは、ストレージ管理装置1、ストレージ装置2から構成される。

20

【0023】

ストレージ管理装置1は、ストレージ装置2を管理する装置である。ストレージ管理装置1は、ストレージ装置2を管理するために、パトロール処理実行手段1a、およびパトロール流量規制手段1bを有している。

【0024】

パトロール処理実行手段1aは、ストレージ装置2の記憶領域2aが正常に稼働していることを確認するパトロール処理を実行する。ここで、記憶領域2aは、ストレージ装置2においてデータを記憶する領域である。このパトロール処理は、具体的には、パトロール処理実行手段1aが、ストレージ装置2の記憶領域2aに記憶されているデータを一旦読み出した後、この読み出したデータを、そのデータが記憶されていた領域に対して書き込み、これを記憶領域2aを区分した領域ごとに順次繰り返すことで行われる。

30

【0025】

パトロール流量規制手段1bは、パトロール処理実行手段1aによって実行されているパトロール処理の速さであるパトロール流量を規制する。このパトロール流量の規制は、具体的には、例えば、パトロール流量規制手段1bが、パトロール処理実行手段1aによって実行されるパトロール処理のパトロール流量を監視し、パトロール流量が速過ぎる場合には、パトロール処理実行手段1aによるパトロール処理の実行を待機させることによ

40

【0026】

このようなストレージ管理装置1によれば、パトロール処理実行手段1aにより、ストレージ装置2の記憶領域2aが正常に稼働していることが確認される。パトロール流量規制手段1bにより、パトロール処理実行手段1aで実行されるパトロール処理の速さが規制される。

【0027】

これによって、パトロール処理の実行に基づく負荷による、ストレージ装置へのアクセスに対する影響を抑えつつ、ストレージシステムの信頼性を高めることができる。

以下、本実施の形態について図面を参照して詳細に説明する。

50

【 0 0 2 8 】

図2は、本実施の形態のシステム構成を示す図である。図2に示すストレージシステムは、データをネットワークで接続された複数のストレージ装置で分散して記憶する。本実施の形態に係るストレージシステムでは、ストレージノード100、200、300、400、コントロールノード500、アクセスノード600および管理ノード30がネットワーク10を介して相互に接続されている。また、端末装置21、22、23が、ネットワーク20を介してアクセスノード600に接続されている。

【 0 0 2 9 】

ストレージノード100、200、300、400には、それぞれストレージ装置110、210、310、410が接続されている。ストレージノード100、200、300、400は、接続されたストレージ装置110、210、310、410に格納されたデータを管理し、管理しているデータをネットワーク10経由でアクセスノード600に提供する。また、ストレージノード100、200、300、400は、データに冗長性をもたせて管理している。すなわち、同一内容のデータが、少なくとも2つのストレージノードで管理されている。

【 0 0 3 0 】

ストレージ装置110には、ハードディスク装置(HDD)111、112、113、114が実装されている。ストレージ装置210には、HDD211、212、213、214が実装されている。ストレージ装置310には、HDD311、312、313、314が実装されている。ストレージ装置410には、HDD411、412、413、414が実装されている。ストレージ装置110、210、310、410は、内蔵する複数のHDDを用いたRAIDシステムである。本実施の形態では、ストレージ装置110、210、310、410は、RAID5のディスク管理サービスを提供する。

【 0 0 3 1 】

なお、ストレージ装置110、210、310、410は、RAID5に限らず、他のRAIDを用いて構成してもよく、RAID以外の技術を用いてディスクアレイを構成してもよい。さらに、ストレージ装置110、210、310、410は、ディスクアレイを用いずにハードディスク単体で構成してもよく、その他の記憶装置を用いてもよい。

【 0 0 3 2 】

ストレージノード100、200、300、400は、接続されたストレージ装置110、210、310、410に格納されたデータを管理し、管理しているデータをネットワーク10経由で端末装置21、22、23に提供する。また、ストレージノード100、200、300、400は、冗長性を有するデータを管理している。すなわち、同一のデータが、少なくとも2つのストレージノードで管理されている。

【 0 0 3 3 】

さらに、ストレージノード100、200、300、400は、パトロール処理において、二重化されたデータの整合性をチェックする。なお、ストレージノード100、200、300、400は個々の判断に基づいてデータの整合性をチェックしてもよいし、外部からの指示によりデータの整合性をチェックしてもよい。本実施の形態では、各ストレージノード100、200、300、400が、自発的に二重化されたデータの整合性をチェックする。

【 0 0 3 4 】

パトロール処理における二重化されたデータの整合性のチェックでは、二重化されたそれぞれのデータを保持するストレージノード同士が互いに通信し合い、冗長性のあるデータの整合性がチェックされる。その際、二重化されたデータを保持する一方のストレージノードで管理されているデータで不具合が検出されれば、他方のストレージノードの対応するデータを用いてデータの復旧が行われる。

【 0 0 3 5 】

コントロールノード500は、ストレージノード100、200、300、400を管理する。具体的には、コントロールノード500は、データの配置状況を示す論理ポリュ

10

20

30

40

50

ームを保持している。コントロールノード500は、ストレージノード100, 200, 300, 400からデータの管理に関する情報を取得し、必要に応じて論理ボリュームを更新する。また、コントロールノード500は、論理ボリュームが更新されると、その影響を受けるストレージノードに対して更新内容を通知する。論理ボリュームに関しては、後述する。

【0036】

アクセスノード600は、端末装置21, 22, 23に対して、ストレージノード100, 200, 300, 400が管理するデータを利用した情報処理のサービスを提供する。すなわち、アクセスノード600は、端末装置21, 22, 23からの要求に回答して所定のプログラムを実行し、必要に応じてストレージノード100, 200, 300, 400にアクセスする。ここで、アクセスノード600は、コントロールノード500から論理ボリュームを取得し、取得した論理ボリュームに基づいてアクセスすべきストレージノードを特定する。

10

【0037】

管理ノード30は、分散ストレージシステムの管理者が操作する端末装置である。分散ストレージシステムの管理者は、管理ノード30を操作して、ストレージノード100, 200, 300, 400、コントロールノード500およびアクセスノード600にアクセスし、運用に必要な各種設定を行うことができる。

【0038】

端末装置21, 22, 23は、ストレージシステムのユーザが操作するコンピュータである。ストレージシステムのユーザは、端末装置21, 22, 23を操作して、ストレージノード100にアクセスする。これに基づいて、ストレージノード100は、ユーザの要求に応じて、ストレージ装置110に記憶されているデータの読み書きを行う。

20

【0039】

次に、ストレージノード100, 200, 300, 400、コントロールノード500、アクセスノード600、端末装置21, 22, 23および管理ノード30のハードウェア構成について説明する。

【0040】

図3は、ストレージノードのハードウェア構成例を示す図である。ストレージノード100は、CPU (Central Processing Unit) 101によって装置全体が制御されている。CPU 101には、バス107を介してRAM (Random Access Memory) 102、HDDインタフェース103、グラフィック処理装置104、入力インタフェース105、および通信インタフェース106が接続されている。

30

【0041】

RAM 102には、CPU 101に実行させるOS (Operating System) のプログラム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、RAM 102には、CPU 101による処理に必要な各種データが格納される。

【0042】

HDDインタフェース103には、ストレージ装置110が接続されている。HDDインタフェース103は、ストレージ装置110に内蔵されたRAIDコントローラ115と通信し、ストレージ装置110に対するデータの入出力を行う。ストレージ装置110内のRAIDコントローラ115は、RAID0~5の機能を有し、複数のHDD111~114をまとめて1台のハードディスクとして管理する。

40

【0043】

グラフィック処理装置104には、モニタ11が接続されている。グラフィック処理装置104は、CPU 101からの命令に従って、画像をモニタ11の画面に表示させる。入力インタフェース105には、キーボード12とマウス13とが接続されている。入力インタフェース105は、キーボード12やマウス13から送られてくる信号を、バス107を介してCPU 101に送信する。

【0044】

50

通信インタフェース 106 は、ネットワーク 10 に接続されている。通信インタフェース 106 は、ネットワーク 10 を介して、他のコンピュータとの間でデータの送受信を行う。

【0045】

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、ストレージノード 100 およびストレージ装置 110 の構成のみを示したが、他のストレージノード 200, 300, 400 や他のストレージ装置 210, 310, 410 も同様のハードウェア構成で実現できる。

【0046】

さらに、コントロールノード 500、アクセスノード 600、および端末装置 21, 22, 23 も、ストレージノード 100 とストレージ装置 110 との組合せと同様のハードウェア構成で実現できる。ただし、コントロールノード 500、アクセスノード 600、および端末装置 21, 22, 23 については、ストレージ装置 110 のような RAID システムではなく、単体の HDD が HDD コントローラに接続されていてもよい。

【0047】

図 2 に示すように、複数のストレージノード 100, 200, 300, 400 がネットワーク 10 に接続され、それぞれのストレージノード 100, 200, 300, 400 は他のストレージノードとの間で通信を行う。この分散ストレージシステムは、端末装置 21, 22, 23 に対して、仮想的なボリューム（以下、論理ボリュームと呼ぶ）として機能する。

【0048】

図 4 は、論理ボリュームのデータ構造を示す図である。論理ボリューム 700 には、「L VOL - A」という識別子（論理ボリューム識別子）が付与されている。また、ネットワーク経由で接続された 4 台のストレージ装置 110, 210, 310, 410 には、それぞれのストレージ装置を管理するストレージノードの識別のためにそれぞれ「SN - A」、「SN - B」、「SN - C」、「SN - D」というノード識別子が付与されている。

【0049】

各ストレージノード 100, 200, 300, 400 が有するストレージ装置 110, 210, 310, 410 それぞれにおいて RAID 5 の論理ディスクが構成されている。この論理ディスクは、5 つのスライスに分割され、個々のストレージノード内で管理されている。

【0050】

図 4 の例では、ストレージ装置 110 内の記憶領域は、5 つのスライス 121 ~ 125 に分けられている。ストレージ装置 210 内の記憶領域は、5 つのスライス 221 ~ 225 に分けられている。ストレージ装置 310 内の記憶領域は、5 つのスライス 321 ~ 325 に分けられている。ストレージ装置 410 内の記憶領域は、5 つのスライス 421 ~ 425 に分けられている。

【0051】

なお、論理ボリューム 700 は、セグメント 710, 720, 730, 740 という単位で構成される。セグメント 710, 720, 730, 740 の記憶容量は、ストレージ装置 110, 210, 310, 410 における管理単位であるスライスの記憶容量と同じである。例えば、スライスの記憶容量が 1 ギガバイトとするとセグメントの記憶容量も 1 ギガバイトである。論理ボリューム 700 の記憶容量はセグメント 1 つ当りの記憶容量の整数倍である。セグメントの記憶容量が 1 ギガバイトならば、論理ボリューム 700 の記憶容量は 4 ギガバイトといったものになる。

【0052】

セグメント 710, 720, 730, 740 は、それぞれプライマリスライス 711, 721, 731, 741 とセカンダリスライス 712, 722, 732, 742 との組から構成される。同じセグメントに属するスライスは別々のストレージノードに属する。個々のスライスを管理する領域には論理ボリューム識別子やセグメント情報や同じセグメン

10

20

30

40

50

トを構成するスライス情報の他にフラグがあり、そのフラグにはプライマリあるいはセカンダリなどを表す値が格納される。

【 0 0 5 3 】

図 4 の例では、スライスの識別子を、「 P 」または「 S 」のアルファベットと数字との組合せで示している。「 P 」はプライマリスライスであることを示している。「 S 」はセカンダリスライスであることを示している。アルファベットに続く数字は、何番目のセグメントに属するのかを表している。例えば、1 番目のセグメント 7 1 0 のプライマリスライスが「 P 1 」で示され、セカンダリスライスが「 S 1 」で示される。

【 0 0 5 4 】

このような構造の論理ボリューム 7 0 0 の各プライマリスライスおよびセカンダリスライスは、ストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 内のいずれかのスライスに対応付けられる。例えば、セグメント 7 1 0 のプライマリスライス 7 1 1 は、ストレージ装置 4 1 0 のスライス 4 2 4 に対応付けられ、セカンダリスライス 7 1 2 は、ストレージ装置 2 1 0 のスライス 2 2 2 に対応付けられている。

【 0 0 5 5 】

そして、各ストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 では、自己のスライスに対応するプライマリスライスまたはセカンダリスライスのデータを記憶する。

次に、ストレージノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 のモジュール構成について説明する。

【 0 0 5 6 】

図 5 は、ストレージノードの機能を示すブロック図である。なお、図 5 ではストレージノード 1 0 0 のモジュール構成を示しているが、他のストレージノード 2 0 0 , 3 0 0 , 4 0 0 もストレージノード 1 0 0 と同様のモジュール構成によって実現できる。

【 0 0 5 7 】

ストレージノード 1 0 0 は、ストレージ装置 1 1 0 を管理する装置である。ストレージノード 1 0 0 は、ストレージ装置 1 1 0 を管理するために、パトロール処理実行部 1 3 1、パトロール流量規制部 1 3 2、管理情報記憶部 1 3 3、パトロール性能測定部 1 3 4、および時刻情報取得部 1 3 5 を有している。同様に、ストレージノード 2 0 0 は、ストレージ装置 2 1 0 を管理する装置である。ストレージノード 2 0 0 は、ストレージ装置 2 1 0 を管理するために、ストレージノード 1 0 0 と同様の機能を有している。

【 0 0 5 8 】

ここで、ストレージ装置 1 1 0 は、ストレージ装置 2 1 0 と二重化されているものとする。すなわち、ストレージ装置 1 1 0 に記憶されているデータには、ストレージ装置 2 1 0 に記憶されている同一の内容の冗長データが存在するものとする。このようにして、ストレージ装置 1 1 0 のデータはストレージ装置 2 1 0 の冗長データによりバックアップされている。

【 0 0 5 9 】

パトロール処理実行部 1 3 1 は、ストレージ装置 1 1 0 の記憶領域 1 1 0 a が正常に稼働していることを確認するパトロール処理を実行する。具体的には、パトロール処理実行部 1 3 1 は、パトロール処理として、ストレージ装置 1 1 0 の記憶領域 1 1 0 a に記憶されているデータを読み出し、データが正常か否かを判断する処理を実行する。ここで、記憶領域 1 1 0 a は、ストレージ装置 1 1 0 においてデータを記憶する領域である。

【 0 0 6 0 】

パトロール処理実行部 1 3 1 は、さらに、ストレージ装置 1 1 0 の記憶領域 1 1 0 a に記憶されているデータの冗長構成が正しく保持されていることを確認する冗長パトロール処理を実行する。具体的には、パトロール処理実行部 1 3 1 は、管理情報記憶部 1 3 3 に記憶されている管理情報を参照して、ストレージ装置 2 1 0 に記憶されている冗長データから、ストレージ装置 1 1 0 に記憶されているデータと同一の内容の冗長データを特定する。次に、パトロール処理実行部 1 3 1 は、ストレージ装置 1 1 0 に記憶されているデータと、ストレージ装置 2 1 0 に記憶されている上記の特定した冗長データとを読み出す。

10

20

30

40

50

次に、パトロール処理実行部 131 は、読み出したデータと冗長データとを比較することによって、比較したデータと冗長データとが一致することを確認する。

【0061】

ここで、パトロール処理実行部 131 は、パトロール処理および冗長パトロール処理を、ストレージ装置 110 の記憶領域 110a が測定単位領域からさらに複数に区分された領域である実行単位領域ごとにパトロールする。

【0062】

また、パトロール処理実行部 131 は、パトロール処理として、ストレージ装置 110 との間で二重化されたストレージ装置 210 に記憶されているデータと同一の内容の冗長データをストレージ装置 110 内のデータの記憶領域 110a に対して、上書きで書き込む処理を実行する。このとき、パトロール処理実行部 131 は、管理情報記憶部 133 に記憶されている管理情報を参照して、ストレージ装置 110 に記憶されているデータと同一の内容の冗長データを特定する。次に、パトロール処理実行部 131 は、ストレージ装置 110 に記憶されているデータと同一の内容であって、ストレージ装置 210 に記憶されている、上記の特定した冗長データを読み出す。次に、パトロール処理実行部 131 は、読み出した冗長データをストレージ装置 110 内のデータの記憶領域 110a に対して、上書きで書き込む処理を実行する。

【0063】

パトロール流量規制部 132 は、パトロール処理実行部 131 によって実行されているパトロール処理および冗長パトロール処理の速さであるパトロール流量を規制する。具体的には、パトロール流量規制部 132 は、パトロール性能測定部 134 によって測定されたパトロール性能に基づいてパトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないようにパトロール流量を規制する。このとき、パトロール流量規制部 132 は、ストレージ装置 110 の記憶領域をパトロール処理が実行される単位に区分した領域である、実行単位領域ごとにパトロール処理のパトロール流量を規制する。

【0064】

また、パトロール流量規制部 132 は、時間帯に応じたパトロール流量の規制を行うこともできる。この場合、パトロール流量規制部 132 は、時刻情報取得部 135 によって取得された時刻情報が示す現在の時刻に基づいて、パトロール性能測定部 134 が保持するパトロール性能から現在の時刻が属する時間帯に対応するパトロール性能を選択する。次に、パトロール流量規制部 132 は、選択した当該パトロール性能に基づいてパトロール流量設定値を決定する。次に、パトロール流量規制部 132 は、決定したパトロール流量設定値を超えないようにパトロール流量を規制する。

【0065】

管理情報記憶部 133 は、ストレージ装置 110 と、当該ストレージ装置 110 に記憶されているデータと同一内容の冗長データを記憶しているストレージ装置 210 とを対応付ける管理情報を記憶する。ストレージノード 100 は、この管理情報に基づいて記憶しているデータの冗長データおよび冗長データが記憶されているストレージ装置を特定する。

【0066】

パトロール性能測定部 134 は、ストレージ装置 110 における読み出し処理と書き込み処理との処理性能を示すパトロール性能を測定する。ここで、パトロール性能測定部 134 は、ストレージ装置 110 の記憶領域 110a が複数に区分された領域である測定単位領域ごとに、パトロール性能を測定する。

【0067】

また、パトロール性能測定部 134 は、時刻情報取得部 135 によって取得された時刻情報が示す時刻が属する時間帯ごとに、測定したパトロール性能を複数種類保持する。

時刻情報取得部 135 は、現在の時刻を示す時刻情報を取得する。この時刻情報は、ストレージノード 100 の CPU 101 が有するタイマによって生成されるが、これに限らず、ネットワーク 10 に接続されたネットワーク時計、または他のノードから取得しても

10

20

30

40

50

よい。

【 0 0 6 8 】

このようなストレージノード 1 0 0 によれば、パトロール処理実行部 1 3 1 により、ストレージ装置 1 1 0 の記憶領域が正常に稼動していることが確認される。このとき、パトロール流量規制部 1 3 2 により、パトロール処理実行部 1 3 1 によって実行されるパトロール処理の速さが規制される。この結果、パトロール処理で生じるストレージシステムへの負荷をコントロールすることによって、ストレージシステムに対するパトロール処理の影響を最小限に留めることができる。

【 0 0 6 9 】

次に、本実施の形態のパトロール性能の測定について説明する。

10

図 6 は、パトロール性能の測定の様子を示す図である。以下、図 6 に従って、本実施の形態のストレージノード 1 0 0 によるパトロール性能の取得について説明する。

【 0 0 7 0 】

本実施の形態におけるパトロール処理の目的は、ストレージ装置 1 1 0 の記憶領域が正常に稼動していること、すなわち、ストレージ装置 1 1 0 の記憶領域に対して正常にアクセスできることを確認することである。そのため、本実施の形態のストレージノードは、パトロール処理において、ストレージ装置 1 1 0 内の記憶領域全体にわたってデータを読み出し、さらにそのデータを書き込む処理を実行する。

【 0 0 7 1 】

また、このパトロール処理がストレージシステムに与える負荷をコントロールするために、そのコントロールの基準とするべく、各ストレージ装置におけるパトロール処理の処理性能を示すパトロール性能を測定する。各ストレージノードでは、ストレージノードの起動時または定期的に、ストレージ装置全体またはストレージ装置の記憶領域のうちの特定の領域について、パトロール性能を測定する。また、各ストレージノードは、測定により得られたパトロール性能を、それぞれのストレージノードが有する RAM (例えば、RAM 1 0 2 など) に格納する。

20

【 0 0 7 2 】

まず、ストレージノード 1 0 0 によって管理される HDD 1 1 1 およびストレージノード 2 0 0 によって管理されるストレージ装置 2 1 0 を構成する HDD 2 1 1 の、それぞれの記憶領域全体のパトロール性能である、全体パトロール性能 1 1 1 p , 2 1 1 p を測定する場合について説明する。この場合、ストレージノード 1 0 0 は、例えば、ストレージシステムの起動時に、HDD 1 1 1 の全体パトロール性能 1 1 1 p を測定する。ストレージノード 2 0 0 は、同様に、HDD 2 1 1 の全体パトロール性能 2 1 1 p を測定する。ストレージノード 1 0 0 , 2 0 0 は、全体パトロール性能 1 1 1 p , 2 1 1 p のように、それぞれ HDD 1 1 1 , 2 1 1 全体について一つのパトロール性能を保持する。そして、ストレージノード 1 0 0 , 2 0 0 は、それぞれの HDD 1 1 1 , 2 1 1 に対するパトロール処理の実行時において、この全体パトロール性能 1 1 1 p , 2 1 1 p をパトロール流量の調整に用いる。

30

【 0 0 7 3 】

次に、ストレージ装置 1 1 0 を構成する HDD 1 1 2 の記憶領域を複数に分割した、それぞれの領域ごとのパトロール性能である領域別パトロール性能 1 1 2 p を測定する場合について説明する。この場合、ストレージノード 1 0 0 は、例えば、ストレージシステムの起動時に、HDD 1 1 2 の領域別パトロール性能 1 1 2 p を測定する。ストレージノード 1 0 0 は、HDD 1 1 2 について、領域別パトロール性能 1 1 2 p のように、複数に分割されたその領域ごとにパトロール性能を保持する。そして、ストレージノード 1 0 0 は、HDD 1 1 2 における各領域に対するパトロール処理の実行時において、この領域別パトロール性能 1 1 2 p を、パトロール流量の調整に用いる。

40

【 0 0 7 4 】

次に、ストレージ装置 1 1 0 を構成する HDD 1 1 3 の全体の時間帯ごとのパトロール性能である、時間帯別パトロール性能 1 1 3 p を測定する場合について説明する。この場

50

合、ストレージノード100は、HDD113全体についてのパトロール性能を、時間帯別パトロール性能113pのように、時間帯ごとに保持する。そして、ストレージノード100は、HDD113に対するパトロール処理の実行時の時間帯に合わせて、この時間帯別パトロール性能113pを、パトロール流量の調整に用いる。

【0075】

次に、本実施の形態のストレージノード100, 200, 300, 400で用いられるパトロール性能テーブルについて説明する。図7は、ストレージ装置ごとのパトロール性能の測定結果を格納したパトロール性能テーブルのデータ構造を示す図である。図8は、分割された領域ごとのストレージ装置のパトロール性能の測定結果を格納したパトロール性能テーブルのデータ構造を示す図である。図9は、時間帯ごとのストレージ装置のパトロール性能の測定結果を格納したパトロール性能テーブルのデータ構造を示す図である。

10

【0076】

図7～図9に示すパトロール性能テーブルは、パトロール処理を行う際の処理を進める速さであるパトロール流量を調整するためのパトロール性能を示す情報である。パトロール情報は、ストレージ装置110, 210, 310, 410が備える各HDDのデータ読み書きの性能特性を示す情報である。図6に示すように取得されたパトロール情報は、パトロール性能テーブル191, 192, 193としてテーブル化されて、ストレージノード100のRAM102およびストレージノード200, 300, 400のRAM(図示省略)に格納されている。

【0077】

20

図7に示すパトロール性能テーブル191は、ストレージ装置ごとのパトロール性能である、全体パトロール性能(図6の全体パトロール性能111pを参照)の測定結果を格納するテーブルである。このパトロール性能テーブル191には、装置191a、シーケンシャル・アクセスの場合の読み出し処理191eおよび書き込み処理191f、ならびにランダム・アクセスの場合の読み出し処理191gおよび書き込み処理191hの欄が設けられており、各欄の横方向に並べられた情報同士が互いに関連づけられている。

【0078】

装置191aの欄には、ストレージ装置110, 210, 310, および410に対して、それぞれが識別可能になるように割り当てられたストレージ装置の名称が格納される。

30

【0079】

シーケンシャル・アクセスの場合の読み出し処理191eおよび書き込み処理191fの欄には、それぞれ、各ストレージ装置におけるシーケンシャル・アクセスによる読み出し処理および書き込み処理の際の全体パトロール性能を示す情報が格納される。シーケンシャル・アクセスによる読み出し処理および書き込み処理の際のパトロール性能は、データ転送速度、すなわち1秒間当りのデータ転送量(MB/s: MegaByte per second)によって示される。MB/sは、大きい値の方が高評価となる。

【0080】

ランダム・アクセスの場合の読み出し処理191gおよび書き込み処理191hの欄には、それぞれ、各ストレージ装置におけるランダム・アクセスによる読み出し処理および書き込み処理の際の全体パトロール性能を示す情報が格納される。ランダム・アクセスによる読み出し処理および書き込み処理のパトロール性能は、処理能力、すなわち、1秒間当りの書き込み処理および読み出し処理の実行回数(IOPS: Input Output Per Second)によって示される。IOPSは、大きい値の方が高評価となる。

40

【0081】

図8に示すパトロール性能テーブル192は、ストレージ装置における分割された領域ごとのパトロール性能である、領域別パトロール性能(図6の領域別パトロール性能112pを参照)の測定結果を格納するテーブルである。このパトロール性能テーブル192には、装置192a、領域192b、シーケンシャル・アクセスの場合の読み出し処理192eおよび書き込み処理192f、ならびにランダム・アクセスの場合の読み出し処理

50

192gおよび書き込み処理192hの欄が設けられており、各欄の横方向に並べられた情報同士が互いに関連づけられている。

【0082】

装置192aの欄には、パトロール性能テーブル191と同様、ストレージ装置110, 210, 310, および410に対して、それぞれが識別可能になるように割り当てられたストレージ装置の名称が格納される。

【0083】

領域192bの欄には、各ストレージ装置における分割された領域の範囲を特定可能な情報が格納される。例えば、図8のパトロール性能テーブル192の一番上の行のように、領域192bが「0~100[block]」である場合には、その行には、ストレージ装置の領域0~100[block]の領域別パトロール性能が格納されることを示す。ここで、blockは、各HDDにおけるデータの管理単位であり、ここでは、1block=1MBを用いるが、必ずしもこれに限られない。

10

【0084】

シーケンシャル・アクセスの場合の読み出し処理192eおよび書き込み処理192fの欄には、それぞれ、各ストレージ装置におけるシーケンシャル・アクセスによる読み出し処理および書き込み処理の際の領域別パトロール性能を示す情報が格納される。ランダム・アクセスの場合の読み出し処理192gおよび書き込み処理192hの欄には、それぞれ、各ストレージ装置におけるランダム・アクセスによる読み出し処理および書き込み処理の際の領域別パトロール性能を示す情報が格納される。

20

【0085】

図9に示すパトロール性能テーブル193は、各ストレージ装置における時間帯ごとのパトロール性能である、時間帯別パトロール性能(図6の時間帯別パトロール性能113pを参照)の測定結果を格納するテーブルである。このパトロール性能テーブル193には、装置193a、対象時間193c、測定日193d、シーケンシャル・アクセスの場合の読み出し処理193eおよび書き込み処理193f、ならびにランダム・アクセスの場合の読み出し処理193gおよび書き込み処理193hの欄が設けられており、各欄の横方向に並べられた情報同士が互いに関連づけられている。

【0086】

装置193aの欄には、パトロール性能テーブル191と同様、ストレージ装置110, 210, 310, および410に対して、それぞれが識別可能になるように割り当てられたストレージ装置の名称が格納される。

30

【0087】

対象時間193cの欄には、各ストレージ装置の時間帯別パトロール性能が測定された時間帯を特定可能な情報が格納される。例えば、図9のパトロール性能テーブル193の一番上の行のように、対象時間193cが「0~3時」である場合には、その行には、0時から3時に測定されたパトロール性能が格納されることを示す。

【0088】

測定日193dの欄には、各ストレージ装置の時間帯別パトロール性能が測定された日付を特定可能な情報が格納される。例えば、図9のパトロール性能テーブル193の一番上の行のように、測定日193dが「10/1」である場合には、その行の時間帯別パトロール性能は、10月1日に測定されたものであることを示す。

40

【0089】

シーケンシャル・アクセスの場合の読み出し処理193eおよび書き込み処理193fの欄には、それぞれ、各ストレージ装置におけるシーケンシャル・アクセスによる読み出し処理および書き込み処理の際の時間帯別パトロール性能を示す情報が格納される。ランダム・アクセスの場合の読み出し処理193gおよび書き込み処理193hの欄には、それぞれ、各ストレージ装置におけるランダム・アクセスによる読み出し処理および書き込み処理の際の時間帯別パトロール性能を示す情報が格納される。

【0090】

50

ストレージノード 100, 200, 300, 400 は、図 11 に示すパトロール処理の実行時、および図 13 に示す冗長パトロール処理の実行時に、対応付けられているストレージ装置 110, 210, 310, 410 からパトロール性能を収集し、それぞれが備える RAM のパトロール性能テーブル 191, 192, 193 に格納する。

【0091】

また、ストレージ装置 110, 210, 310, 410 のパトロール性能として、例えば、応答時間など、その他の性能を用いてもよい。

ここで、パトロール性能とパトロール流量との関係について説明する。図 7 ~ 図 9 に示すパトロール性能テーブル 191, 192, 193 に格納されているパトロール性能に基づいて、例えば、以下の(例 1) ~ (例 3) に示すようなポリシーに従い、ストレージノード 100 により、パトロール流量が決定される。ここでは、“Seq-R” は、シーケンシャル・アクセスによる読み出しであり、“Seq-W” は、シーケンシャル・アクセスによる書き込みである。また、“Ran-R” は、ランダム・アクセスによる読み出しであり、“Ran-W” は、ランダム・アクセスによる書き込みである。

10

【0092】

(例 1) 「アクセスに影響を与えない範囲で可能な限りパトロールを行う」場合
パトロール処理の実行開始前の 1 分間の平均アクセスが、

Seq-R 10[MB/s]

Seq-W 0[MB/s]

Ran-R 100[IOPS]

Ran-W 0[IOPS]

20

であり、ストレージ装置 110 の最大性能が、

Seq-R 50[MB/s]

Seq-W 50[MB/s]

Ran-R 500[IOPS]

Ran-W 500[IOPS]

である場合、シーケンシャル・アクセスの読み出しで、ストレージ装置 110 の最大性能の 20% を占めており、ランダム・アクセスの読み出しで、最大性能の 20% を占めている。このため、パトロール処理に対しては、ストレージ装置 110 の最大性能の 60% を割り当てることができる。

30

【0093】

ここで、パトロール処理では、シーケンシャル・アクセスの読み出しおよび書き込みが平行して行われる。このため、パトロール流量は、それぞれ、シーケンシャル・アクセスの読み出しおよび書き込みに対して割り当て可能な 60% の二分の一となり、15 [MB/s] となる。

【0094】

(例 2) 「ランダム・アクセスが性能の 50% を超えているときには、アクセスに影響を与えない範囲で可能な限りパトロールを行い、50% 以下の場合にはパトロール処理に性能の 20% を費やすようにする」場合

パトロール処理の実行開始前の 1 分間の平均アクセスが、

Seq-R 20[MB/s]

Seq-W 0[MB/s]

Ran-R 100[IOPS]

Ran-W 100[IOPS]

40

であり、性能値が、

Seq-R 50[MB/s]

Seq-W 50[MB/s]

Ran-R 500[IOPS]

Ran-W 500[IOPS]

である場合、ランダム・アクセスは、読み出しおよび書き込みにおいて、それぞれパト

50

ール性能の20%を占めている。このランダム・アクセスの読み出しおよび書き込みは、合わせてストレージ装置110の最大性能の40%を占めている。すなわち、この時点のランダム・アクセスは、ストレージ装置110の最大性能の50%以下である。

【0095】

ここで、パトロール処理では、シーケンシャル・アクセスの読み出しおよび書き込みが平行して行われる。このため、ディスク性能の20%を割り当てると、10 [MB/s]となることから、パトロール流量は、その二分の一の5 [MB/s]となる。

【0096】

(例3)「アクセス状況に関係なく、常に1 [MB/s]でパトロール処理を行う」場合
パトロール流量は、常に1 [MB/s]となる。

10

次に、ストレージ装置の記憶領域において、パトロール性能が測定される分割単位およびパトロール処理が実行される分割単位の関係について説明する。図10は、パトロール処理の実行時におけるストレージ装置の領域の分割を示す図である。ここでは、ストレージ装置110の有する記憶領域を一体として、記憶領域116として説明する。

【0097】

図10に示す、ストレージ装置110の記憶領域116は、パトロール性能の測定の際の分割単位である複数の測定単位116a, 116b, 116c, 116d, …, 116nに分割されている。図8に示す領域別パトロール性能は、この測定単位116a, 116b, 116c, 116d, …, 116nごとに測定される。

20

【0098】

また、記憶領域116は、さらに、パトロール処理が実行される際の分割単位である複数の実行単位116a1, 116a2, …, 116anに分割されている。後述する図11に示すパトロール処理および図14に示す冗長パトロール処理は、この実行単位116a1, 116a2, …, 116anごとに処理が進められ、パトロール流量の規制が行われる。

【0099】

次に、本実施の形態のストレージノードにおいて実行される処理の手順について説明する。まず、本実施の形態のストレージノードにおいて実行されるパトロール処理について説明する。図11は、パトロール処理の手順を示すフローチャートである。

30

【0100】

本実施の形態のストレージノード100では、例えば、ストレージシステムの起動時、ストレージノードまたはストレージ装置の起動時、所定の時間の経過時などに、このパトロール処理の実行が開始される。なお、パトロール処理は、ストレージ装置について記憶領域に異常がないかチェックするためにパトロールを行う処理である。

【0101】

以下に、図11に従って、ストレージノード100が、ストレージ装置110についてパトロールを行う場合について説明する。ここでは、ストレージノード100の起動時において、ストレージ装置110の全領域についてパトロールが行われる場合に基づいて説明する。ストレージノード100の起動時において、ストレージノード100は、ストレージ装置110の全領域が正常に読み書き可能であるか否かを診断するパトロール処理を実行するために、パトロール処理を呼び出す。

40

【0102】

[ステップS11]ストレージノード100は、パトロール流量を設定する。このパトロール流量(図6参照)は、上記のように、パトロール処理が行われる速さを調整するための条件である。ここで、ストレージノード100は、パトロール性能として、全体パトロール性能(例えば、図6に示す全体パトロール性能111p, 211p)、領域別パトロール性能(例えば、図6に示す領域別パトロール性能112p)、および時間帯別パトロール性能(例えば、図6に示す時間帯別パトロール性能113p)を用いることができる。ストレージノード100が各ストレージ装置に対していずれのパトロール性能を用い

50

るかについては、各ストレージ装置のハードウェア特性、使用状況に合わせて予め管理者などにより設定される。

【 0 1 0 3 】

例えば、記憶領域によってパトロール性能が大きく異なるストレージ装置については、各領域ごとのパトロール性能を示す領域別パトロール性能を用いる。また、時間帯によって使用状況に偏りがあるストレージ装置については、時間帯ごとのパトロール性能を示す時間帯別パトロール性能を用いる。また、パトロール性能が変化しないストレージ装置については、全体パトロール性能を用いる。ストレージノード 1 0 0 は、これらのパトロール性能を予め測定しておき、パトロール処理の実行時に処理の実行対象のストレージ装置が有する H D D のパトロール性能を取得する。そして、ストレージノード 1 0 0 は、取得したパトロール性能を、そのストレージ装置における各 H D D のパトロール流量として設定する。

10

【 0 1 0 4 】

[ステップ S 1 2] ストレージノード 1 0 0 は、ストレージ装置 1 1 0 について記憶領域の実行単位ごとにパトロールするパトロール単位実行処理を実行する。詳しくは、図 1 2 に示すパトロール単位実行処理において後述する。また、ストレージノード 1 0 0 は、後述する指定時間の起算点とするために、パトロール単位実行処理の実行開始時に、時刻情報を取得する。

【 0 1 0 5 】

[ステップ S 1 3] ストレージノード 1 0 0 は、ステップ S 1 1 において設定されたパトロール流量以下でパトロールが実行されているか否かを判定する。ストレージノード 1 0 0 は、設定されたパトロール流量を超えていれば、ステップ S 1 4 に処理を進める一方、設定されたパトロール流量以下であれば、ステップ S 1 5 に処理を進める。

20

【 0 1 0 6 】

ステップ S 1 3 において、具体的には、ストレージノード 1 0 0 は、実行単位 (図 1 0 参照) ごとのパトロール処理の時間間隔が、設定された指定時間以下であるか否かを判定する。

【 0 1 0 7 】

ここで、指定時間は、パトロール処理がストレージシステムに与える負荷を調整するべく、実行単位ごとのパトロール処理の時間間隔が一定時間以上になるように規制するための値である。この指定時間は、H D D のパトロール性能が高ければパトロール流量を増加 (パトロール処理の速さの制限を緩和) させ、パトロール性能が低ければパトロール流量を低下 (パトロール処理の速さの制限を強化) させるものを用いることができる。

30

【 0 1 0 8 】

本実施の形態では、指定時間を、以下の関係式を用いて算出する。

$$\text{指定時間[s]} = \text{実行単位[MB]} / \text{パトロール性能[MB/s]} \quad (1)$$

具体的には、例えば、ある H D D の全体パトロール性能が 1 0 [MB/s] の場合において、この H D D の実行単位が 1 [MB] のときは、指定時間は、 $1 \text{ [MB]} / 1 0 \text{ [MB/s]} = 0 . 1 \text{ [s]}$ となる。

【 0 1 0 9 】

本実施の形態のストレージノード 1 0 0 は、以上のように指定時間をこのように算出する。そして、ストレージノード 1 0 0 は、直近のパトロール単位実行処理 (S 1 2) の開始時点から起算して指定時間が経過していれば、「設定流量以下である」と判定する一方、前回のパトロール処理の開始時点から指定時間が経過していなければ、「設定流量を超えた」と判定する。

40

【 0 1 1 0 】

[ステップ S 1 4] ストレージノード 1 0 0 は、直近のパトロール単位実行処理の開始時点から起算して指定時間が経過するまでステップ S 1 4 において処理を待機する。ストレージノード 1 0 0 は、指定時間が経過した後、ステップ S 1 5 に処理を進める。

【 0 1 1 1 】

50

具体的には、ストレージノード100は、前回のパトロール処理が終了してから一定時間（例えば、0.1秒）経過するまで、ステップS14において処理を待機する。

【ステップS15】ストレージノード100は、ステップS12のパトロール単位実行処理が、ストレージ装置110のすべての実行単位について実行されたか否かを判定する。ストレージノード100は、すべての実行単位について実行されていない場合は、ステップS12に処理を進める一方、すべての実行単位について実行されていれば、処理を終了する。

【0112】

次に、上記のパトロール処理（図11参照）のステップS12において実行されるパトロール単位実行処理の詳細について説明する。図12は、パトロール単位実行処理の手順を示すフローチャートである。ここでは、図11に示すパトロール処理と同様、ストレージノード100の起動時において、ストレージ装置110の全領域についてパトロールが行われる場合に基づいて説明する。

10

【0113】

【ステップS21】ストレージノード100は、ストレージ装置110の記憶領域における現在のパトロール対象の実行単位に格納されているデータを読み出す。

【ステップS22】ストレージノード100は、ステップS21における読み出しに成功したか否かを判定する。ストレージノード100は、読み出しに成功していれば、ステップS23に処理を進める一方、読み出しに失敗していれば、ステップS24に処理を進める。

20

【0114】

【ステップS23】ストレージノード100は、ステップS21において読み出したデータを、ステップS21においてデータを読み出した記憶領域に書き込む。その後、ストレージノード100は、パトロール単位実行処理を終了してパトロール処理に復帰する。

【0115】

【ステップS24】ストレージノード100は、ステップS21において読み出したデータと対応する冗長データを、冗長データが記憶されているストレージ装置から読み出して、ステップS21においてデータを読み出した記憶領域に書き込む。その後、ストレージノード100は、パトロール単位実行処理を終了してパトロール処理に復帰する。

30

【0116】

以上のようにして、本実施の形態のストレージノード100は、パトロール処理を実行することにより、ストレージ装置110の記憶領域が正常に稼動していること、すなわち、ストレージ装置110の記憶領域に正常に読み書きできることを確認する。また、パトロール処理を自律的に調整して実行することで、パトロール処理に基づくストレージシステムの負荷を規制する。この結果、パトロール処理の実行に基づく負荷による、ストレージ装置へのアクセスに対する影響を抑えつつ、ストレージシステムの信頼性を高めることができる。

【0117】

次に、本実施の形態のストレージノード100で実行される冗長パトロールについて説明する。本実施の形態のストレージノード100では、上記のパトロール処理に加えて、冗長パトロール処理が実行される。この冗長パトロール処理は、ストレージ装置110に記憶されているデータをバックアップするためにそのデータと同一の内容である冗長データがストレージ装置110と二重化されている他のストレージ装置（例えば、ストレージ装置210）に記憶されている場合には、ストレージ装置110に記憶されているデータと、そのデータに対応する冗長データの整合性について確認する処理である。

40

【0118】

本実施の形態のストレージノード100では、冗長パトロール処理をストレージノード100の起動時および定期的に行う。このとき、冗長パトロール処理は、ストレージシステムに対する負荷を調整するために、冗長パトロール性能に基づいて実行される。

【0119】

50

次に、本実施の形態の冗長パトロール性能の測定について説明する。

図13は、冗長パトロール性能の測定の様子を示す図である。以下、図13に従って、本実施の形態のストレージノード100による冗長パトロール性能の取得について説明する。

【0120】

本実施の形態における冗長パトロール処理の目的は、ストレージ装置110の記憶領域に記憶されているデータが、対応する冗長データと整合し、冗長構成が維持されていること、すなわち、ストレージ装置110の記憶されているデータが二重化されているストレージ装置によって正常にバックアップされていることを確認することである。そのため、本実施の形態のストレージノードは、冗長パトロール処理によって、ストレージ装置110内の記憶領域全体にわたってデータを読み出し、さらにそのデータを書き込む処理を実行する。

10

【0121】

また、パトロール処理と同様、この冗長パトロール処理がストレージシステムに与える負荷をコントロールするために、そのコントロールの基準とするべく、冗長パトロール処理の処理性能を示すパトロール性能である冗長パトロール性能を測定する。このとき、冗長構成をとるストレージ装置を管理するストレージノード間で、協調して冗長パトロール処理を進める必要があるため、パトロール性能としてノード間の冗長構成も考慮に入れた値である冗長パトロール性能を測定し、それを冗長パトロール処理における処理の速さ(冗長パトロール流量)の調整に用いる。

20

【0122】

各ストレージノードでは、ストレージノードの起動時または定期的に、ストレージ装置全体またはストレージ装置の記憶領域のうち特定の領域について、二重化された他方のストレージ装置を管理するストレージノードと協調して、冗長パトロール性能を測定する。また、各ストレージノードは、測定により得られた冗長パトロール性能を、それぞれのストレージノードが有するRAM(例えば、RAM102など)に格納する。

【0123】

図13は、ストレージノード100が管理するストレージ装置110のHDD111とストレージノード200が管理するストレージ装置210のHDD211とが二重化(RAID1)構成となっているときの冗長パトロール性能511pを測定する場合の測定の様子について図示したものである。この場合、ストレージノード100,200は、冗長パトロール性能511pのように、HDD111,211の二つのHDDの組ごとに一つの冗長パトロール性能値を保持する。

30

【0124】

次に、本実施の形態のストレージノードにおいて実行される冗長パトロール処理について説明する。冗長パトロール処理は、ストレージ装置に記憶されているデータと、冗長構成のストレージ装置に記憶されている冗長データとの整合性をチェックするためにパトロールを行う処理である。図14は、冗長パトロール処理の手順を示すフローチャートである。

【0125】

以下に、図14に従って、ストレージノード100が、ストレージノード200とともに、ストレージ装置110および冗長構成のストレージ装置210について冗長パトロールを行う場合について説明する。この冗長パトロール処理は、二重化されたストレージ装置のうち一方のストレージ装置(例えば、ストレージ装置110)を管理するストレージノード(例えば、ストレージノード100)が主導して実行するが、他方のストレージ装置(例えば、ストレージ装置210)を管理するストレージノード(例えば、ストレージノード200)と協調して処理が進められる。ここでは、ストレージノード100の起動時におけるパトロール処理が実行された後、ストレージ装置110の全領域、およびこれに冗長構成として関係付けられているストレージ装置210の全領域についてパトロールが行われる場合に基づいて説明する。

40

50

【 0 1 2 6 】

ストレージノード 1 0 0 の起動時においてパトロール処理が実行された後、ストレージノード 1 0 0 は、ストレージ装置 1 1 0 の全領域が、ストレージ装置 1 1 0 のデータと冗長構成として関係付けられたストレージ装置 2 1 0 のデータとの間で、整合性を有しているか否かを診断する冗長パトロール処理を実行するために、冗長パトロール処理を呼び出す。

【 0 1 2 7 】

[ステップ S 3 1] ストレージノード 1 0 0 は、パトロール処理を実行する二重化された HDD の組に対応する冗長パトロール流量 (図 1 3 参照) を設定する。

[ステップ S 3 2] ストレージノード 1 0 0 は、ストレージ装置 1 1 0 について記憶領域の実行単位ごとにパトロールするとともに、ストレージ装置 2 1 0 を管理するストレージノードがストレージ装置 2 1 0 をパトロールしたパトロール結果を受信して、両者の整合性をチェックする冗長パトロール単位実行処理を実行する。詳しくは、図 1 5 に示す冗長パトロール単位実行処理において後述する。また、ストレージノード 1 0 0 は、後述する指定時間の起算点とするために、冗長パトロール単位実行処理の実行開始時に、時刻情報を取得する。

【 0 1 2 8 】

[ステップ S 3 3] ストレージノード 1 0 0 は、パトロール処理のステップ S 1 3 と同様に、ステップ S 3 1 において設定されたパトロール流量以下で冗長パトロールが実行されているか否かを判定する。ストレージノード 1 0 0 は、設定されたパトロール流量を超えていれば、ステップ S 3 4 に処理を進める一方、設定されたパトロール流量以下であれば、ステップ S 3 5 に処理を進める。

【 0 1 2 9 】

[ステップ S 3 4] ストレージノード 1 0 0 は、直近の冗長パトロール単位実行処理 (ステップ S 3 2) の開始時点から起算して指定時間が経過するまでステップ S 3 4 において処理を待機する。ストレージノード 1 0 0 は、指定時間が経過した後、ステップ S 3 5 に処理を進める。

【 0 1 3 0 】

[ステップ S 3 5] ストレージノード 1 0 0 は、ステップ S 3 2 の冗長パトロール単位実行処理が、ストレージ装置 1 1 0 のすべての実行単位について実行されたか否かを判定する。ストレージノード 1 0 0 は、すべての実行単位について実行されていなければ、ステップ S 3 2 に処理を進める一方、すべての実行単位について実行されていれば、処理を終了する。

【 0 1 3 1 】

次に、上記の冗長パトロール処理 (図 1 4 参照) のステップ S 3 2 において実行される冗長パトロール単位実行処理の詳細について説明する。図 1 5 は、冗長パトロール単位実行処理の手順を示すシーケンス図である。ここでは、ストレージノード 1 0 0 の起動時におけるパトロール処理が実行された後、ストレージ装置 1 1 0 の全領域、およびこれに冗長構成として関係付けられているストレージ装置 2 1 0 の全領域についてパトロールが行われる場合に基づいて説明する。

【 0 1 3 2 】

[ステップ S 4 1] ストレージノード 1 0 0 は、冗長構成としてストレージ装置 1 1 0 と関係付けられたストレージ装置 2 1 0 を管理するストレージノード 2 0 0 に対して、ストレージ装置 2 1 0 に対する冗長パトロール処理の開始を要求する冗長パトロール要求を送信する。このとき、ストレージノード 1 0 0 は、併せて、ストレージノード 2 0 0 に対して、冗長パトロール処理の対象であるストレージ装置 2 1 0 の実行単位を特定するための情報を送信する。

【 0 1 3 3 】

[ステップ S 4 2] ストレージノード 1 0 0 は、ストレージ装置 1 1 0 の記憶領域における現在のパトロール対象の実行単位に格納されているデータを読み出す。

[ステップS43]ストレージノード100は、ストレージノード200から送信された、ストレージ装置210に記憶されていた冗長データを受信する。

【0134】

[ステップS44]ストレージノード100は、ステップS42において読み出したデータとステップS43で受信したデータの整合性をチェックする。

[ステップS45]ストレージノード100は、ステップS44において行った整合性のチェック結果を、ストレージノード200に送信する。

【0135】

[ステップS46]ストレージノード100は、ステップS44において行った整合性のチェック結果を、ストレージノード100のRAM102に保存する。その後、ストレージノード100は、冗長パトロール単位実行処理を終了して冗長パトロール処理に復帰する。

【0136】

[ステップS51]ストレージノード200は、ステップS41においてストレージノード100から送信された冗長パトロール要求を受信する。このとき、ストレージノード200は、併せて送信された情報から、冗長パトロール処理の対象であるストレージ装置210の実行単位を特定する。

【0137】

[ステップS52]ストレージノード200は、ストレージ装置210の記憶領域におけるステップS51において特定された実行単位に格納されているデータを読み出す。このデータは、ストレージ装置110に記憶されている、冗長パトロールの対象であるデータの冗長データである。

【0138】

[ステップS53]ストレージノード200は、ステップS52においてストレージ装置210から読み出した冗長データを、ストレージノード100に対して送信する。

[ステップS54]ストレージノード200は、ステップS45においてストレージノード100から送信された整合性のチェック結果を受信する。

【0139】

[ステップS55]ストレージノード200は、ステップS54において受信した整合性のチェック結果を、ストレージノード200のRAM(図示省略)に保存する。その後、ストレージノード200は、処理を終了する。

【0140】

以上のようにして、本実施の形態のストレージノードは、冗長パトロール処理を実行することにより、ストレージノード間で冗長構成をとっている場合に、各ノードのデータが正しく冗長構成を保持しているか否かを診断することができる。

【0141】

以上、本発明のストレージ管理装置、ストレージ管理プログラムおよびストレージシステムを、図示の実施の形態に基づいて説明したが、上記については単に本発明の原理を示すものである。本発明は上記に示し、説明した正確な構成および応用例に限定されるものではなく、さらに、多数の変形、変更が当業者にとって可能であり、対応するすべての変形例および均等物は、添付の請求項およびその均等物による本発明の範囲とみなされ、各部の構成は同様の機能を有する任意の構成のものに置換することができる。また、本発明に他の任意の構成物や行程が付加されてもよい。また、本発明は前述した実施の形態のうちの任意の2以上の構成(特徴)を組み合わせたものであってもよい。

【0142】

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、ストレージノード100, 200, 300, 400、管理ノード30、端末装置21, 22, 23、コントロールノード500およびアクセスノード600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。

10

20

30

40

50

【 0 1 4 3 】

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。コンピュータで読み取り可能な記録媒体には、磁気記録装置、光ディスク、光磁気記録媒体、半導体メモリなどがある。磁気記録装置には、HDD、フレキシブルディスク(FD)、磁気テープ(MT)などがある。光ディスクには、DVD(Digital Versatile Disc)、DVD-RAM、CD-ROM(Compact Disc - Read Only Memory)、CD-R(Recordable)/RW(ReWritable)などがある。光磁気記録媒体には、MO(Magneto - Optical disk)などがある。

【 0 1 4 4 】

上記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、CD-ROM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に格納しておき、ネットワークを通じて、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。

10

【 0 1 4 5 】

上記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若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

20

【 0 1 4 6 】

以上の実施の形態に関し、さらに以下の付記を開示する。

(付記1) データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムにおけるストレージ装置を管理するストレージ管理装置において、

前記ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段と、

前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段と、

を有することを特徴とするストレージ管理装置。

30

【 0 1 4 7 】

(付記2) 前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、

前記パトロール処理実行手段は、さらに、前記パトロール処理において、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記ストレージ装置に記憶されている前記データと、前記二重化ストレージ装置に記憶されている前記データと同一内容の前記冗長データとを読み出し、読み出した前記データと前記冗長データとを比較することによって、比較した前記データと前記冗長データとが一致することを確認することを特徴とする付記1記載のストレージ管理装置。

40

【 0 1 4 8 】

(付記3) 前記ストレージ装置における読み出し処理と書き込み処理との処理性能を示すパトロール性能を測定するパトロール性能測定手段を有し、

前記パトロール流量規制手段は、前記パトロール性能測定手段によって測定された前記パトロール性能に基づいて前記パトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないように前記パトロール流量を規制することを特徴とする付記1記載のストレージ管理装置。

【 0 1 4 9 】

(付記4) 前記パトロール性能測定手段は、前記ストレージ装置の記憶領域が複数に

50

区分された領域である測定単位領域ごとに、前記パトロール性能を測定し、

前記パトロール処理実行手段は、前記パトロール処理を、前記ストレージ装置の記憶領域が前記測定単位領域からさらに複数に区分された領域である実行単位領域ごとに実行し、

前記パトロール流量規制手段は、前記実行単位領域ごとに前記パトロール処理の前記パトロール流量を規制することを特徴とする付記 1 記載のストレージ管理装置。

【 0 1 5 0 】

(付記 5) 現在の時刻を示す時刻情報を取得する時刻情報取得手段を有し、

前記パトロール性能測定手段は、前記時刻情報取得手段によって取得された前記時刻情報が示す時刻が属する時間帯ごとに、測定した前記パトロール性能を複数種類保持し、

前記パトロール流量規制手段は、前記時刻情報取得手段によって取得された前記時刻情報が示す前記現在の時刻に基づいて、前記パトロール性能測定手段が保持する前記パトロール性能から前記現在の時刻が属する時間帯に対応する前記パトロール性能を選択し、選択した当該パトロール性能に基づいて前記パトロール流量設定値を決定し、決定した当該パトロール流量設定値を超えないように前記パトロール流量を規制することを特徴とする付記 3 記載のストレージ管理装置。

【 0 1 5 1 】

(付記 6) 前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、

前記パトロール処理実行手段は、前記パトロール処理として、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記二重化ストレージ装置に記憶されている前記データと同一内容の前記冗長データとを読み出し、読み出した前記冗長データを前記ストレージ装置内の前記データの記憶領域に対して上書きで書き込む処理を実行することを特徴とする付記 1 記載のストレージ管理装置。

【 0 1 5 2 】

(付記 7) コンピュータに、データをネットワークで接続された複数のストレージ装置に分散して記憶するストレージシステムにおけるストレージ装置の管理処理を実行させるストレージ管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の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段、

前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段、

として機能させることを特徴とするストレージ管理プログラム。

【 0 1 5 3 】

(付記 8) データをネットワークで接続された複数のストレージ装置に分散して記憶し、前記ストレージ装置を管理するストレージ管理装置を有するストレージシステムにおいて、

前記ストレージ装置と、当該ストレージ装置に記憶されている前記データと同一内容の冗長データを記憶している二重化ストレージ装置とを対応付ける管理情報を記憶する管理情報記憶手段を有し、

前記ストレージ管理装置は、

前記ストレージ装置の記憶領域が正常に稼動していることを確認するパトロール処理を実行するパトロール処理実行手段と、

前記パトロール処理実行手段によって実行されている前記パトロール処理の速さであるパトロール流量を規制するパトロール流量規制手段と、

を有し、

前記パトロール処理実行手段は、さらに、前記パトロール処理において、前記管理情報記憶手段に記憶されている前記管理情報を参照して、前記ストレージ装置に記憶されてい

10

20

30

40

50

る前記データと、前記二重化ストレージ装置に記憶されている前記データと同一内容の前記冗長データを読み出し、読み出した前記データと前記冗長データとを比較することによって、比較した前記データと前記冗長データとが一致することを確認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】

【0154】

【図1】本実施の形態の概要を示す図である。

【図2】本実施の形態のシステム構成を示す図である。

【図3】ストレージノードのハードウェア構成を示す図である。

【図4】論理ボリュームのデータ構造を示す図である。

【図5】ストレージノードの機能を示すブロック図である。

【図6】パトロール性能の測定の様子を示す図である。

【図7】ストレージ装置ごとのパトロール性能の測定結果を格納したパトロール性能テーブルのデータ構造を示す図である。

【図8】分割された領域ごとのストレージ装置のパトロール性能の測定結果を格納したパトロール性能テーブルのデータ構造を示す図である。

【図9】時間帯ごとのストレージ装置のパトロール性能の測定結果を格納したパトロール性能テーブルのデータ構造を示す図である。

【図10】パトロール処理の実行時におけるストレージ装置の領域の分割を示す図である。

。

【図11】パトロール処理の手順を示すフローチャートである。

【図12】パトロール単位実行処理の手順を示すフローチャートである。

【図13】冗長パトロール性能の測定の様子を示す図である。

【図14】冗長パトロール処理の手順を示すフローチャートである。

【図15】冗長パトロール単位実行処理の手順を示すシーケンス図である。

【符号の説明】

【0155】

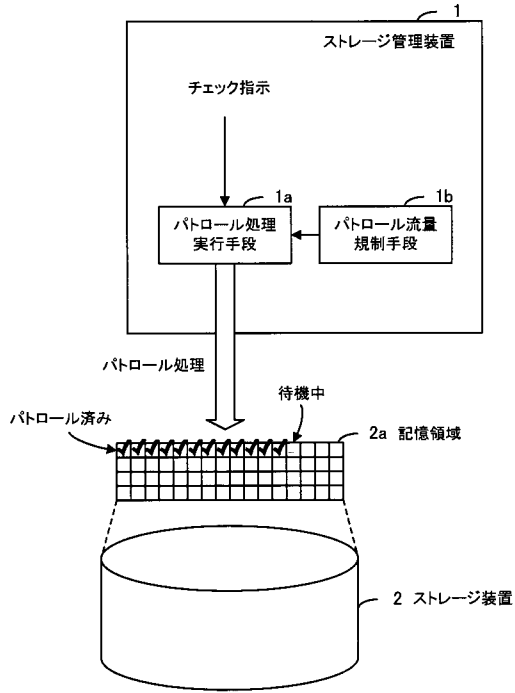
- 1 ストレージ管理装置
 - 1 a パトロール処理実行手段
 - 1 b パトロール流量規制手段
- 2 ストレージ装置
 - 2 a 記憶領域

10

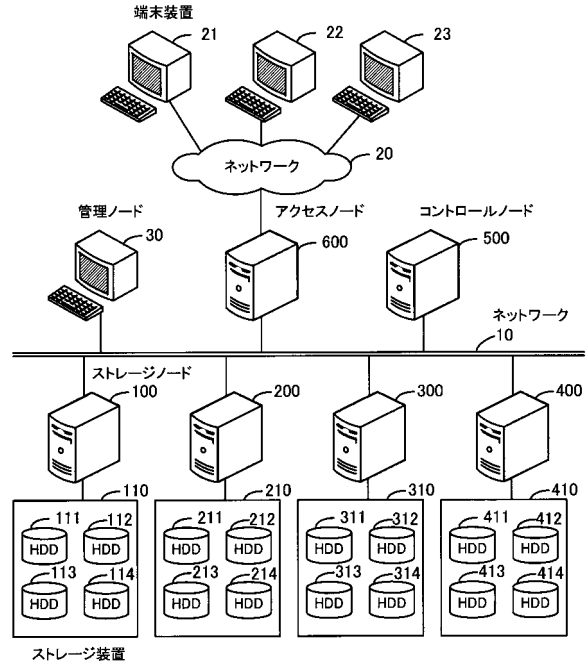
20

30

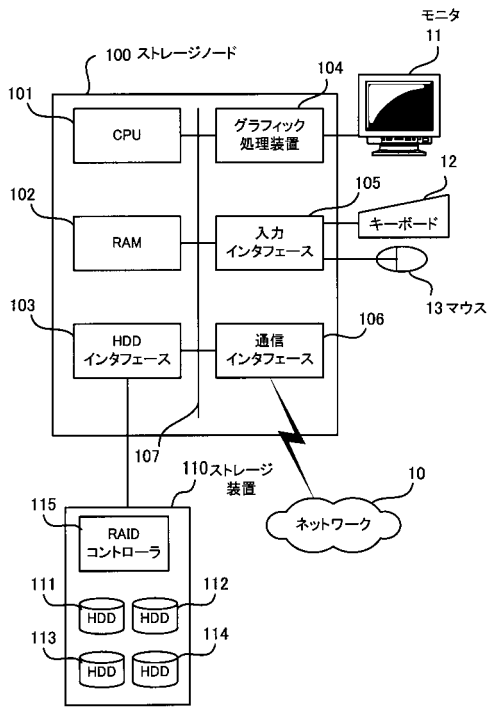
【図1】



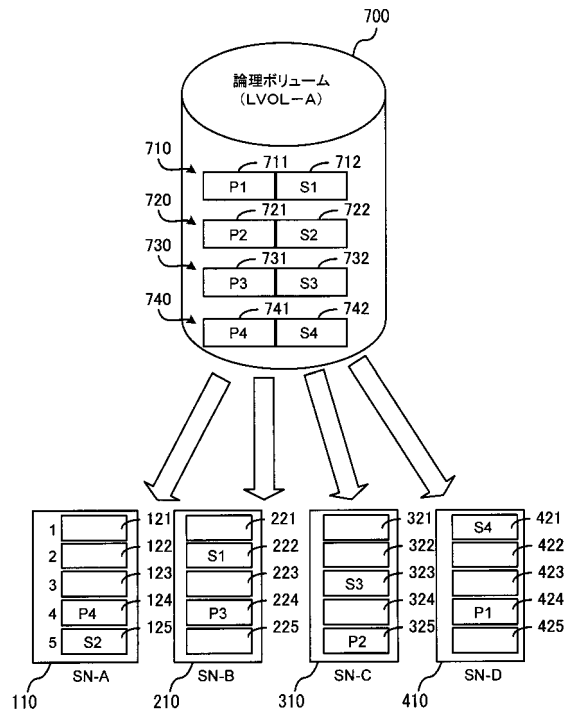
【図2】



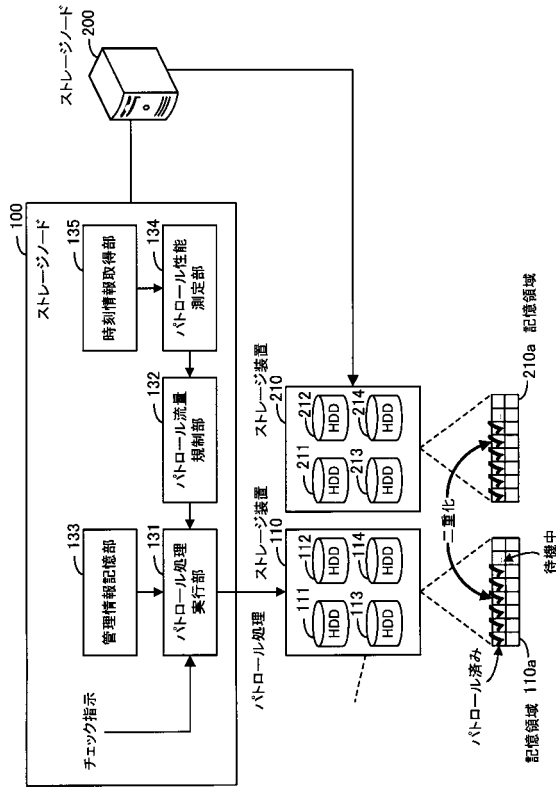
【図3】



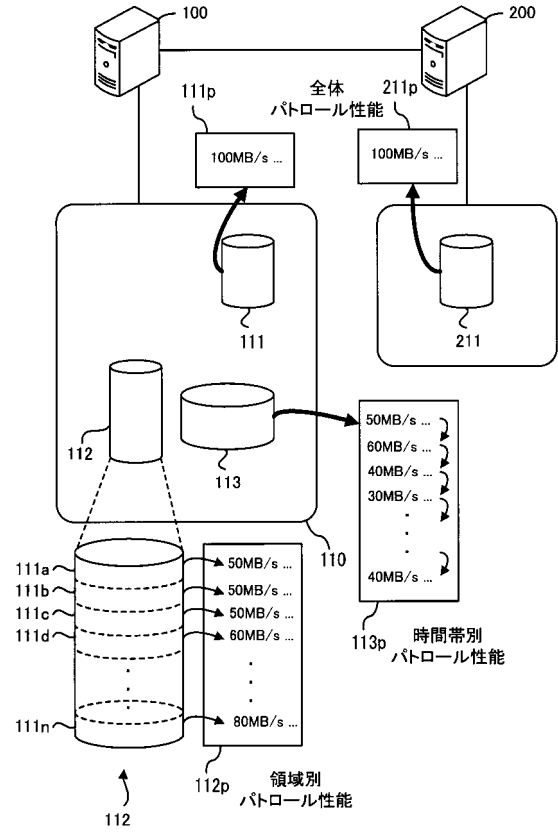
【図4】



【図5】



【図6】



【図7】

パトロール性能テーブル			
装置	シーケンシャル		ランダム
	読み出し処理	書き込み処理	書き込み処理
ストレージ装置A	50 MB/s	50 MB/s	500 IOPS
ストレージ装置B	30 MB/s	10 MB/s	10 IOPS
...
191a	191e	191f	191h

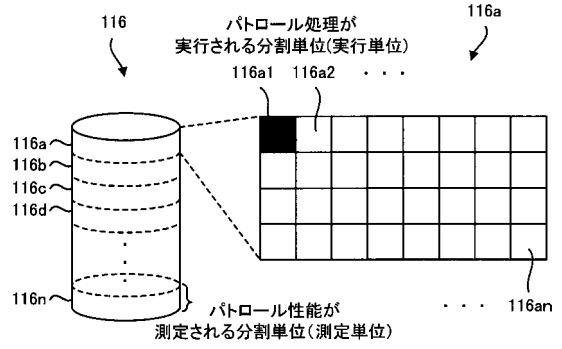
【図8】

パトロール性能テーブル				
装置	領域	シーケンシャル		ランダム
		読み出し処理	書き込み処理	書き込み処理
ストレージ装置A	0~100[block]	50 MB/s	50 MB/s	500 IOPS
	100~200[block]	30 MB/s	10 MB/s	10 IOPS
	200~1000[block]	30 MB/s	10 MB/s	10 IOPS
...
192a	192b	192e	192f	192h

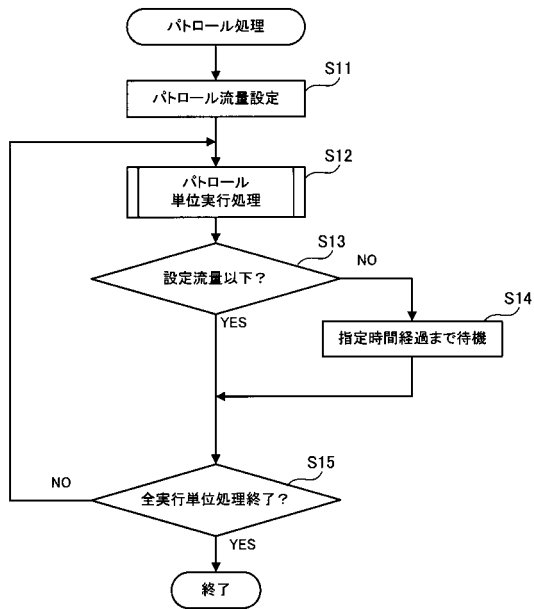
【図9】

バトロール性能テーブル		ランダム	
装置	対象時間	シーケンシャル	
		読み出し処理	書き込み処理
ストレージ装置A	0~3時	50 MB/s	500 IOPS
	3~6時	50 MB/s	500 IOPS
	6~9時	50 MB/s	200 IOPS
	9~12時	30 MB/s	100 IOPS
	12~15時	10 MB/s	10 IOPS
	15~19時	30 MB/s	100 IOPS
	19~21時	30 MB/s	100 IOPS
	21~24時	30 MB/s	200 IOPS
...

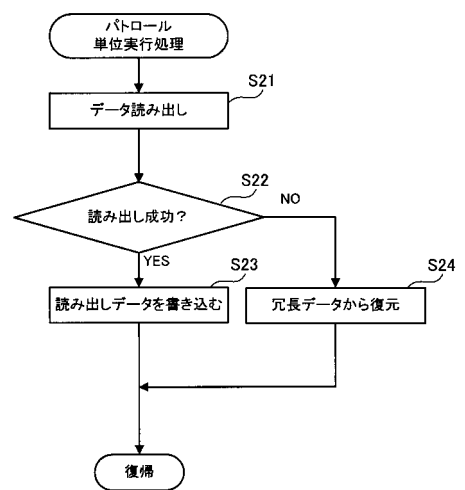
【図10】



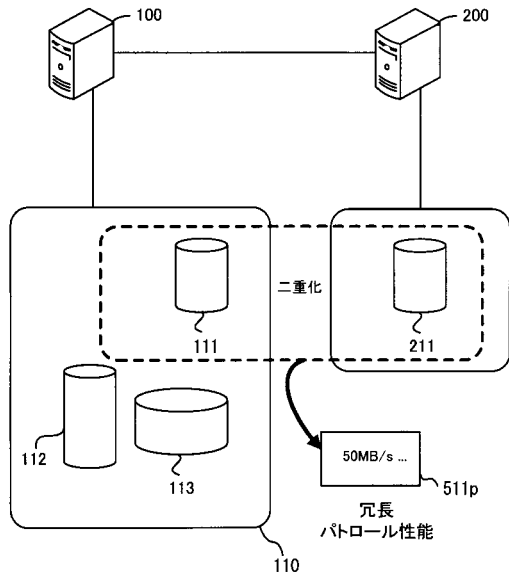
【図11】



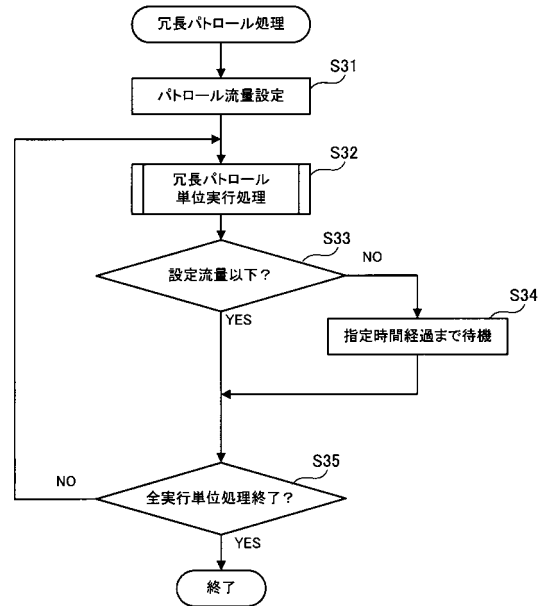
【図12】



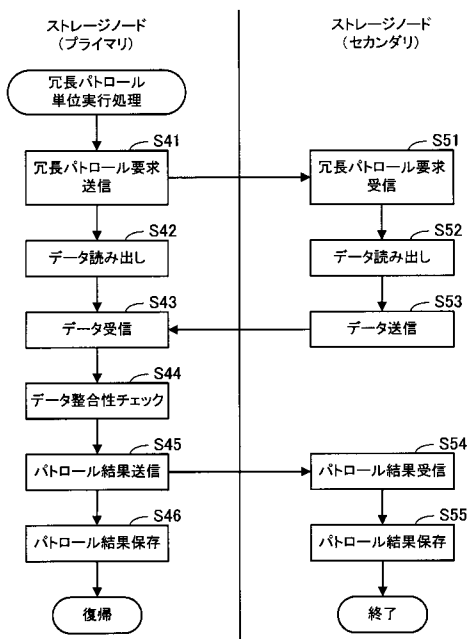
【図13】



【図14】



【図15】



フロントページの続き

- (72)発明者 土屋 芳浩
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
- (72)発明者 丸山 哲太郎
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
- (72)発明者 武 理一郎
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 木村 貴俊

- (56)参考文献 特開平07-210325(JP,A)
特開平09-062461(JP,A)
特開平07-104947(JP,A)
特開2005-157933(JP,A)
国際公開第2004/104845(WO,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08
G06F 12/00 - 12/16
G06F 13/10 - 13/14