

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5099128号
(P5099128)

(45) 発行日 平成24年12月12日(2012.12.12)

(24) 登録日 平成24年10月5日(2012.10.5)

(51) Int.Cl.		F 1			
G06F 12/00	12/00	(2006.01)	G06F	12/00	531D
G06F 3/06	3/06	(2006.01)	G06F	3/06	301P
			G06F	12/00	545A

請求項の数 18 (全 27 頁)

(21) 出願番号	特願2009-512810 (P2009-512810)	(73) 特許権者	000005223
(86) (22) 出願日	平成19年4月20日 (2007.4.20)		富士通株式会社
(86) 国際出願番号	PCT/JP2007/058633		神奈川県川崎市中原区上小田中4丁目1番1号
(87) 国際公開番号	W02008/136075	(74) 代理人	100092152
(87) 国際公開日	平成20年11月13日 (2008.11.13)		弁理士 服部 毅巖
審査請求日	平成21年7月3日 (2009.7.3)	(72) 発明者	荻原 一隆
前置審査			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	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	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

最終頁に続く

(54) 【発明の名称】 ストレージ管理プログラム、ストレージ管理装置およびストレージ管理方法

(57) 【特許請求の範囲】

【請求項1】

同一内容の複数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ノードに分散して配置する分散ストレージシステムのデータ配置状況を管理するストレージ管理プログラムにおいて、

コンピュータを、

前記同一内容の複数のデータから、アクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとを示しており前記主データおよび前記副データそれぞれの配置先の前記ストレージノードを登録した管理情報を記憶する管理情報記憶手段、

前記ストレージノードの負荷情報を継続的に収集する負荷情報収集手段、

前記管理情報記憶手段が記憶する前記管理情報と前記負荷情報収集手段が収集した前記負荷情報とに基づいて、前記主データの配置先のストレージノードと前記副データの配置先のストレージノードとの間で各配置先のストレージノードの負荷の差が所定の許容量を超える、同一内容の前記主データと前記副データとの組を交換対象として検出する交換対象検出手段、

前記交換対象検出手段が検出した1組のデータの間で前記主データと前記副データとの役割を交換するように、前記管理情報記憶手段が記憶する前記管理情報を更新する管理情報更新手段、

前記複数のストレージノードのうちアクセス先のストレージノードを前記管理情報に基

づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが前記副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると前記管理情報の参照要求を送信し当該参照要求に応じて取得した前記管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、前記参照要求を受信すると、前記アクセスノードに更新後の前記管理情報を提供する手段、

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。

【請求項 2】

前記交換対象検出手段は、前記ストレージノードのうち負荷最大ノードと負荷最小ノードとの間で負荷の差が前記所定の許容量を超えるか否かを判断し、前記所定の許容量を超える場合、前記負荷最大ノードに配置された前記主データと前記負荷最小ノードに配置された前記副データとを交換対象の候補とすることを特徴とする請求項 1 記載のストレージ管理プログラム。

10

【請求項 3】

前記交換対象検出手段は、交換対象の検出処理を所定の時間間隔で継続して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プログラム。

【請求項 4】

前記管理情報では、データの特用に用いる論理ボリュームのアドレス空間が複数の論理セグメントに分割され、前記論理セグメント単位で前記同一内容の複数のデータが管理されており、

前記管理情報更新手段は、前記論理セグメント単位で前記主データと前記副データとの役割を交換する、

20

ことを特徴とする請求項 1 記載のストレージ管理プログラム。

【請求項 5】

前記コンピュータを、更に、前記ストレージノードから当該ストレージノードに配置されたデータの格納位置の再編成を開始する通知を受け付ける通知受付手段として機能させ

、前記交換対象検出手段は、前記通知受付手段が通知を受け付けると、通知元の前記ストレージノードに配置された全ての前記主データを特定し、負荷の大きさに拘わらず、特定した前記主データと当該主データと同一内容の前記副データとの組を交換対象とする、

ことを特徴とする請求項 1 記載のストレージ管理プログラム。

30

【請求項 6】

前記コンピュータを、更に、前記管理情報更新手段が前記管理情報を更新した後、役割が変更されたデータが配置されている前記ストレージノードに対して変更内容を通知する更新内容通知手段として機能させることを特徴とする請求項 1 記載のストレージ管理プログラム。

【請求項 7】

同一内容の複数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ノードに分散して配置する分散ストレージシステムのデータ配置状況を管理するストレージ管理装置において、

前記同一内容の複数のデータから、アクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとを示しており前記主データおよび前記副データそれぞれの配置先の前記ストレージノードを登録した管理情報を記憶する管理情報記憶手段と、

40

前記ストレージノードの負荷情報を継続的に収集する負荷情報収集手段と、

前記管理情報記憶手段が記憶する前記管理情報と前記負荷情報収集手段が収集した前記負荷情報とに基づいて、前記主データの配置先のストレージノードと前記副データの配置先のストレージノードとの間で各配置先のストレージノードの負荷の差が所定の許容量を超える、同一内容の前記主データと前記副データとの組を交換対象として検出する交換対象検出手段と、

前記交換対象検出手段が検出した 1 組のデータの間で前記主データと前記副データとの

50

役割を交換するように、前記管理情報記憶手段が記憶する前記管理情報を更新する管理情報更新手段と、

前記複数のストレージノードのうちアクセス先のストレージノードを前記管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが前記副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると前記管理情報の参照要求を送信し当該参照要求に応じて取得した前記管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、前記参照要求を受信すると、前記アクセスノードに更新後の前記管理情報を提供する手段と、
を有することを特徴とするストレージ管理装置。

【請求項 8】

前記交換対象検出手段は、前記ストレージノードのうち負荷最大ノードと負荷最小ノードとの間で負荷の差が前記所定の許容量を超えるか否か判断し、前記所定の許容量を超える場合、前記負荷最大ノードに配置された前記主データと前記負荷最小ノードに配置された前記副データとを交換対象の候補とすることを特徴とする請求項 7 記載のストレージ管理装置。

【請求項 9】

前記交換対象検出手段は、交換対象の検出処理を所定の時間間隔で継続して実行することを特徴とする請求項 7 記載のストレージ管理装置。

【請求項 10】

前記管理情報では、データの特定に用いる論理ボリュームのアドレス空間が複数の論理セグメントに分割され、前記論理セグメント単位で前記同一内容の複数のデータが管理されており、

前記管理情報更新手段は、前記論理セグメント単位で前記主データと前記副データとの役割を交換する、

ことを特徴とする請求項 7 記載のストレージ管理装置。

【請求項 11】

前記ストレージノードから当該ストレージノードに配置されたデータの格納位置の再編成を開始する通知を受け付ける通知受付手段を更に有し、

前記交換対象検出手段は、前記通知受付手段が通知を受け付けると、通知元の前記ストレージノードに配置された全ての前記主データを特定し、負荷の大きさに拘わらず、特定した前記主データと当該主データと同一内容の前記副データとの組を交換対象とする、

ことを特徴とする請求項 7 記載のストレージ管理装置。

【請求項 12】

前記管理情報更新手段が前記管理情報を更新した後、役割が変更されたデータが配置されている前記ストレージノードに対して変更内容を通知する更新内容通知手段を更に有することを特徴とする請求項 7 記載のストレージ管理装置。

【請求項 13】

同一内容の複数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ノードに分散して配置する分散ストレージシステムのデータ配置状況を管理するコンピュータによるストレージ管理方法において、前記コンピュータが、

前記ストレージノードの負荷情報を継続的に収集し、

前記同一内容の複数のデータからアクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとを示しており前記主データおよび前記副データそれぞれの配置先の前記ストレージノードを登録した管理情報記憶手段が記憶する管理情報と、収集した前記負荷情報とに基づいて、前記主データの配置先のストレージノードと前記副データの配置先のストレージノードとの間で各配置先のストレージノードの負荷の差が所定の許容量を超える、同一内容の前記主データと前記副データとの組を交換対象として検出し、

検出した 1 組のデータの間で前記主データと前記副データとの役割を交換するように、前記管理情報記憶手段が記憶する前記管理情報を更新し、

10

20

30

40

50

前記複数のストレージノードのうちアクセス先のストレージノードを前記管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが前記副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると前記管理情報の参照要求を送信し当該参照要求に応じて取得した前記管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、前記参照要求を受信すると、前記アクセスノードに更新後の前記管理情報を提供する、

ことを特徴とするストレージ管理方法。

【請求項 14】

交換対象の検出処理では、前記ストレージノードのうち負荷最大ノードと負荷最小ノードとの間で負荷の差が前記所定の許容量を超えるか否か判断し、前記所定の許容量を超える場合、前記負荷最大ノードに配置された前記主データと前記負荷最小ノードに配置された前記副データとを交換対象の候補とすることを特徴とする請求項 13 記載のストレージ管理方法。

10

【請求項 15】

交換対象の検出処理は所定の時間間隔で継続して実行することを特徴とする請求項 13 記載のストレージ管理方法。

【請求項 16】

前記管理情報では、データの特定に用いる論理ボリュームのアドレス空間が複数の論理セグメントに分割され、前記論理セグメント単位で前記同一内容の複数のデータが管理されており、

20

前記管理情報の更新処理では、前記論理セグメント単位で前記主データと前記副データとの役割を交換する、

ことを特徴とする請求項 13 記載のストレージ管理方法。

【請求項 17】

交換対象の検出処理では、前記ストレージノードから当該ストレージノードに配置されたデータの格納位置の再編成を開始する通知を受け付けていたときは、通知元の前記ストレージノードに配置された全ての前記主データを特定し、負荷の大きさに拘わらず、特定した前記主データと当該主データと同一内容の前記副データとの組を交換対象とすることを特徴とする請求項 13 記載のストレージ管理方法。

【請求項 18】

30

前記管理情報の更新処理の後、役割が変更されたデータが配置されている前記ストレージノードに対して変更内容を通知することを特徴とする請求項 13 記載のストレージ管理方法。

【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のデータ配置状況を管理するストレージ管理プログラム、ストレージ管理装置およびストレージ管理方法に関し、特に複数のストレージノードに分散してデータを配置する分散ストレージシステムのデータ配置状況を管理するストレージ管理プログラム、ストレージ管理装置およびストレージ管理方法に関する。

40

【背景技術】

【0002】

現在、コンピュータによる情報処理が普及するに伴い、データを蓄積するストレージシステムの高機能化の要求が一層高まっている。例えば、大量のデータを読み書きする処理が高速であること、すなわち、高性能化が求められている。また、ストレージシステムを構成するハードウェアの一部に障害が発生しても蓄積したデータが消失しないこと、すなわち、高信頼化が求められている。

【0003】

このような高性能化・高信頼化の要求に応えるストレージシステムとして、分散ストレージシステムが知られている。分散ストレージシステムでは、ネットワークで接続された

50

複数のストレージノードにデータが分散して配置される。また、分散ストレージシステムでは、同一内容のデータが複数のストレージノードに重複して配置されることが多い。すなわち、複数のストレージノードを用いて、負荷の分散化およびデータの冗長化が行われる。これにより、ストレージシステムの応答性能および信頼性を高めることができる。

【0004】

ところで、分散ストレージシステムの運用時には、データの追加や削除に伴って、データの配置状況が絶えず変化する。ここで、頻繁にアクセス要求があるデータが特定のストレージノードに集中して配置されると、負荷の分散が不十分となり、応答性能が低下する。従って、分散ストレージシステムの応答性能を高く維持するためには、運用時に個々のストレージノードの負荷を考慮してデータの配置を管理する必要がある。

10

【0005】

具体的には、分散ストレージシステムの応答性能を高く維持する管理方法として、以下の方法が知られている。第1の方法は、アクセス要求を受け付けたときに、動的にアクセス先のストレージノードを判定する方法である。すなわち、分散ストレージシステムは、アクセス要求があったデータが複数のストレージノードに重複して配置されている場合、アクセス要求時に最も負荷が小さいストレージノードをアクセス先とする（例えば、特許文献1, 2参照）。一方、第2の方法は、ストレージノード間でアクセス要求の量ができる限り均等になるように、データ自体を自動的に移動する方法である（例えば、特許文献3参照）。これにより、特定のストレージノードに負荷が集中することを防止できる。

【特許文献1】特開昭63-56873号公報

20

【特許文献2】特開平11-161555号公報

【特許文献3】特開2005-276017号公報

【発明の開示】

【発明が解決しようとする課題】

【0006】

しかし、上記特許文献1~3に記載の技術には以下のような問題がある。特許文献1, 2に記載の動的な負荷分散の技術では、アクセス要求元とストレージノードとの間で、アクセス要求毎に必ずアクセス先を判定する処理が実行される。このため、データの配置管理を実現する仕組みが複雑になると共に、アクセス先を判定する処理だけ応答が遅くなる。また、特許文献3に記載のデータの自動移動の技術では、データの移動自体がストレージノードに負荷を与える。このため、データ移動中の応答が遅くなる。また、データの移動を繰り返すと、ディスク装置などのハードウェアへの負担が大きくなり、故障の原因ともなる。

30

【0007】

本発明はこのような点に鑑みてなされたものであり、データの冗長化がなされている分散ストレージシステムにおいて、ストレージノード間の負荷分散を容易に実現するストレージ管理プログラム、ストレージ管理装置およびストレージ管理方法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】

【0008】

40

本発明では、上記課題を解決するために、図1に示すストレージ管理プログラムが提供される。本発明に係るストレージ管理プログラムは、同一内容の複数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ノードに分散して配置する分散ストレージシステムのデータ配置状況を管理するものである。ストレージ管理プログラムを実行するコンピュータ1は、管理情報記憶手段1a、負荷情報収集手段1b、交換対象検出手段1c、管理情報更新手段1dおよび管理情報を提供する手段を有する。管理情報記憶手段1aは、同一内容の複数のデータから、アクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとを示しており主データおよび副データそれぞれの配置先のストレージノードを登録した管理情報を記憶する。負荷情報収集手段1bは、ストレージノード2, 3, 4の負荷情報を継続的に収集する。交換対象検出手段1cは、管理情報

50

記憶手段 1 a が記憶する管理情報と負荷情報収集手段 1 b が収集した負荷情報とに基づいて、主データの配置先のストレージノードと副データの配置先のストレージノードとの間で各配置先のストレージノードの負荷の差が所定の許容量を超える、同一内容の主データと副データとの組を交換対象として検出する。管理情報更新手段 1 d は、交換対象検出手段 1 c が検出した 1 組のデータの間で主データと副データとの役割を交換するように、管理情報記憶手段 1 a が記憶する管理情報を更新する。管理情報を提供する手段は、複数のストレージノードのうちアクセス先のストレージノードを管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると管理情報の参照要求を送信し当該参照要求に応じて取得した管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、参照要求を受信すると、アクセスノードに、更新後の管理情報を提供する。

10

【 0 0 0 9 】

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 1 によれば、負荷情報収集手段 1 b により、ストレージノード 2, 3, 4 の負荷情報が継続的に収集される。交換対象検出手段 1 c により、主データの配置先の負荷と副データの配置先の負荷との差が所定の許容量を超えるような、同一内容の主データと副データとの組が検出される。そして、管理情報更新手段 1 d により、検出された 1 組のデータの間で主データと副データの役割が交換される。すなわち、交換前にバックアップ用であったデータがアクセス用のデータとなり、交換前にアクセス用であったデータがバックアップ用のデータとなる。更に、管理情報を提供する手段により、複数のストレージノードのうちアクセス先のストレージノードを管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると管理情報の参照要求を送信し当該参照要求に応じて取得した管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、参照要求を受信されると、アクセスノードに、更新後の管理情報が提供される。

20

【 0 0 1 0 】

また、上記課題を解決するために、同一内容の複数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ノードに分散して配置する分散ストレージシステムのデータ配置状況を管理するストレージ管理装置において、同一内容の複数のデータから、アクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとを示しており主データおよび副データそれぞれの配置先のストレージノードを登録した管理情報を記憶する管理情報記憶手段と、ストレージノードの負荷情報を継続的に収集する負荷情報収集手段と、管理情報記憶手段が記憶する管理情報と負荷情報収集手段が収集した負荷情報とに基づいて、主データの配置先のストレージノードと副データの配置先のストレージノードとの間で各配置先のストレージノードの負荷の差が所定の許容量を超える、同一内容の主データと副データとの組を交換対象として検出する交換対象検出手段と、交換対象検出手段が検出した 1 組のデータの間で主データと副データとの役割を交換するように、管理情報記憶手段が記憶する管理情報を更新する管理情報更新手段と、複数のストレージノードのうちアクセス先のストレージノードを管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると管理情報の参照要求を送信し当該参照要求に応じて取得した管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、参照要求を受信すると、アクセスノードに、更新後の管理情報を提供する手段と、を有することを特徴とするストレージ管理装置が提供される。

30

40

【 0 0 1 1 】

このようなストレージ管理装置によれば、負荷情報収集手段により、ストレージノードの負荷情報が継続的に収集される。交換対象検出手段により、主データの配置先のストレージノードの負荷と副データの配置先のストレージノードの負荷との差が所定の許容量を

50

超えるような、同一内容の主データと副データとの組が検出される。そして、管理情報更新手段により、検出された1組のデータの間で主データと副データの役割が交換される。すなわち、交換前にバックアップ用であったデータがアクセス用のデータとなり、交換前にアクセス用であったデータがバックアップ用のデータとなる。更に、複数のストレージノードのうちアクセス先のストレージノードを管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると管理情報の参照要求を送信し当該参照要求に応じて取得した管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、参照要求が受信されると、アクセスノードに、更新後の管理情報が提供される。

10

【0012】

また、上記課題を解決するために、同一内容の複数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ノードに分散して配置する分散ストレージシステムのデータ配置状況を管理するコンピュータによるストレージ管理方法において、コンピュータが、ストレージノードの負荷情報を継続的に収集し、同一内容の複数のデータからアクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとを示しており主データおよび副データそれぞれの配置先のストレージノードを登録した管理情報記憶手段が記憶する管理情報と、収集した負荷情報とに基づいて、主データの配置先のストレージノードと副データの配置先のストレージノードとの間で各配置先のストレージノードの負荷の差が所定の許容量を超える、同一内容の主データと副データとの組を交換対象として検出し、検出した1組のデータの間で主データと副データとの役割を交換するように、管理情報記憶手段が記憶する管理情報を更新し、複数のストレージノードのうちアクセス先のストレージノードを管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると管理情報の参照要求を送信し当該参照要求に応じて取得した管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、参照要求を受信すると、アクセスノードに、更新後の管理情報を提供することを特徴とするストレージ管理方法が提供される。

20

【0013】

このようなストレージ管理方法によれば、まずストレージノードの負荷情報が継続的に収集される。次に、主データの配置先のストレージノードの負荷と副データの配置先のストレージノードの負荷との差が所定の許容量を超えるような、同一内容の主データと副データとの組が検出される。そして、検出された1組のデータの間で主データと副データの役割が交換される。すなわち、交換前にバックアップ用であったデータがアクセス用のデータとなり、交換前にアクセス用であったデータがバックアップ用のデータとなる。更に、複数のストレージノードのうちアクセス先のストレージノードを管理情報に基づいて選択し当該アクセス先のストレージノードに格納されたアクセス先のデータが副データであるために当該アクセス先のストレージノードにより当該データへのアクセスが拒否されると管理情報の参照要求を送信し当該参照要求に応じて取得した管理情報に基づいてアクセス先のストレージノードを再選択するアクセスノードから、参照要求が受信されると、アクセスノードに、更新後の管理情報が提供される。

30

40

【発明の効果】**【0014】**

本発明では、アクセス用の主データが配置されたストレージノードとバックアップ用の副データが配置されたストレージノードとの間で負荷が所定の許容値を超えると、設定を変更し、主データと副データの役割を交換することとした。これにより、データを移動することなく、当該データのアクセス先のストレージノードが変更される。このとき、アクセス要求単位でアクセス先が判定されるわけではない。従って、負荷分散の処理自体が原因で応答が遅くなることを防止し、分散ストレージシステムの応答性能を確実に高めるこ

50

とができる。

【0015】

本発明の上記および他の目的、特徴および利点は本発明の例として好ましい実施の形態を表す添付の図面と関連した以下の説明により明らかになるであろう。

【図面の簡単な説明】

【0016】

【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】コントロールノードのハードウェア構成を示す図である。

【図5】論理ボリュームのデータ構造を示す第1の模式図である。

【図6】ストレージノードの機能を示すブロック図である。

【図7】コントロールノードおよびアクセスノードの機能を示すブロック図である。

【図8】スライス情報テーブルのデータ構造例を示す図である。

【図9】ボリューム情報テーブルのデータ構造例を示す図である。

【図10】負荷情報テーブルのデータ構造例を示す図である。

【図11】第1の種別変更処理の手順を示すフローチャートである。

【図12】第2の種別変更処理の手順を示すフローチャートである。

【図13】論理ボリュームのデータ構造を示す第2の模式図である。

【図14】種別一斉変更処理の手順を示すフローチャートである。

【図15】アクセスノードからのデータアクセスの流れを示すシーケンス図である。

【発明を実施するための最良の形態】

【0017】

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

図1は、本実施の形態の概要を示す図である。図1に示す分散ストレージシステムは、コンピュータ1、ストレージノード2, 3, 4、コンピュータ5およびネットワーク6から構成される。コンピュータ1、ストレージノード2, 3, 4およびコンピュータ5は、ネットワーク6に接続されている。

【0018】

コンピュータ1は、ストレージノード2, 3, 4に配置されたデータの配置状況を管理する。コンピュータ1は、管理情報記憶手段1a、負荷情報収集手段1b、交換対象検出手段1cおよび管理情報更新手段1dを有する。

【0019】

管理情報記憶手段1aには、管理情報が格納されている。管理情報は、データの配置状況を管理する情報である。具体的には、管理情報には、データそれぞれの配置先のストレージノードが登録されている。また、同一内容の複数のデータがストレージノード2, 3, 4に分散されて配置されている場合、アクセス要求時にアクセス先として使用する主データとバックアップとして使用する副データとが指定されている。

【0020】

負荷情報収集手段1bは、ストレージノード2, 3, 4の負荷情報を収集する。負荷情報には、例えば、CPU (Central Processing Unit) の負荷、受け付けたアクセス要求の数、ネットワーク使用率などが含まれる。負荷情報収集手段1bは、最新の負荷情報を継続的に収集する。

【0021】

交換対象検出手段1cは、管理情報記憶手段1aに格納された管理情報と、負荷情報収集手段1bによって収集された負荷情報とに基づいて、所定の条件を具備する同一内容の主データと副データとの組を検出する。所定の条件とは、主データの配置先のストレージノードの負荷と、副データの配置先のストレージノードの負荷との差が、所定の許容量を超えることである。交換対象検出手段1cは、上記条件を具備する主データと副データと

10

20

30

40

50

の組が存在するか否かを調べる処理を継続的に実行する。

【 0 0 2 2 】

管理情報更新手段 1 d は、交換対象検出手段 1 c によって上記条件を具備する主データと副データとの組が検出されると、管理情報記憶手段 1 a に格納された管理情報を更新する。具体的には、管理情報更新手段 1 d は、検出された 1 組のデータの間で、主データと副データとの役割を交換する。すなわち、交換前にバックアップ用であったデータがアクセス用のデータとなり、交換前にアクセス用であったデータがバックアップ用のデータとなる。

【 0 0 2 3 】

コンピュータ 5 は、コンピュータ 1 から管理情報を取得し、管理情報に基づいてストレージノード 2 , 3 , 4 に配置されたデータにアクセスする。例えば、利用しようとするデータの主データがストレージノード 2 に配置され副データがストレージノード 3 に配置されているとき、コンピュータ 5 はストレージノード 2 にアクセスする。アクセスによって主データが更新されたときは、ストレージノード間で通信が行われ、更新内容が自動的に副データに反映される。ただし、コンピュータ 5 が主データと副データとを同時に更新するようにしてもよい。

【 0 0 2 4 】

ここで、データ # 1 の主データがストレージノード 2 (管理情報では“ A ”と表記)に副データがストレージ 3 (管理情報では“ B ”と表記)に配置され、データ # 2 の主データがストレージノード 3 に副データがストレージノード 4 (管理情報では“ C ”と表記)に配置され、データ # 3 の主データがストレージノード 2 に副データがストレージノード 4 に配置されているとする。また、ストレージノード 2 , 3 , 4 のうち、ストレージノード 2 の負荷が最も高くストレージノード 4 の負荷が最も低く、両者の負荷の差が所定の許容量を超えているとする。

【 0 0 2 5 】

このとき、コンピュータ 1 は、管理情報を更新して、データ # 3 の主データと副データとの役割を交換する。すなわち、交換後は、データ # 3 の主データがストレージノード 4 に副データがストレージノード 2 に配置されることになる。これによって、コンピュータ 5 はデータ # 3 を利用するとき、ストレージノード 4 にアクセスすることになる。

【 0 0 2 6 】

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 1 によれば、負荷情報収集手段 1 b により、ストレージノード 2 , 3 , 4 の負荷情報が継続的に収集される。交換対象検出手段 1 c により、主データの配置先の負荷と副データの配置先の負荷との差が所定の許容量を超えるような、同一内容の主データと副データとの組が検出される。そして、管理情報更新手段 1 d により、検出された 1 組のデータの間で主データと副データの役割が交換される。すなわち、交換前にバックアップ用であったデータがアクセス用のデータとなり、交換前にアクセス用であったデータがバックアップ用のデータとなる。

【 0 0 2 7 】

これにより、データを移動することなく、当該データのアクセス先のストレージノードが変更される。このとき、アクセス要求単位でアクセス先が判定されるわけではない。従って、負荷分散の処理自体が原因で応答が遅くなることを防止し、分散ストレージシステムの応答性能を確実に高めることができる。

【 0 0 2 8 】

以下、本実施の形態を図面を参照して詳細に説明する。

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 に示す分散ストレージシステムは、同一内容の複数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ノードに分散して配置することで、信頼性と処理性能とを向上させたストレージシステムである。

【 0 0 2 9 】

本実施の形態に係る分散ストレージシステムでは、ストレージノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0、コントロールノード 5 0 0、アクセスノード 6 0 0 および管理ノード 3

10

20

30

40

50

0 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。

【0030】

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。

10

【0031】

ストレージ装置 110 には、複数の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、複数の HDD 211, 212, 213, 214 が実装されている。ストレージ装置 310 には、複数の HDD 311, 312, 313, 314 が実装されている。ストレージ装置 410 には、複数の HDD 411, 412, 413, 414 が実装されている。各ストレージ装置 110, 210, 310, 410 は、内蔵する複数の HDD を用いた RAID (Redundant Array of Independent Disks) システムである。本実施の形態では、各ストレージ装置 110, 210, 310, 410 は、RAID 5 のディスク管理サービスを提供する。

【0032】

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ボリュ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。論理ボリュームについては、後で詳細に説明する。

20

【0033】

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。

30

【0034】

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、運用時に必要な各種設定を行うことができる。

40

【0035】

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、管理ノード 30 および端末装置 21, 22, 23 のハードウェア構成について説明する。

【0036】

図 3 は、ストレージノードのハードウェア構成を示す図である。ストレージノード 100 は、CPU 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。

50

【0037】

RAM102には、CPU101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、RAM102には、CPU101による処理に必要な各種データが格納される。

【0038】

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~RAID5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。

10

【0039】

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、CPU101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CPU101に送信する。

【0040】

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。

20

【0041】

ストレージノード200,300,400も、ストレージノード100と同様のハードウェア構成によって実現できる。

図4は、コントロールノードのハードウェア構成を示す図である。コントロールノード500は、CPU501によって装置全体が制御されている。CPU501には、バス507を介してRAM502、HDD503、グラフィック処理装置504、入力インタフェース505および通信インタフェース506が接続されている。

【0042】

RAM502には、CPU501に実行させるOS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、RAM502には、CPU501による処理に必要な各種データが格納される。HDD503には、OSやアプリケーションのプログラムが格納される。

30

【0043】

グラフィック処理装置504には、モニタ51が接続されている。グラフィック処理装置504は、CPU501からの命令に従って、画像をモニタ51の画面に表示させる。入力インタフェース505には、キーボード52とマウス53とが接続されている。入力インタフェース505は、キーボード52やマウス53から送られてくる信号を、バス507を介してCPU501に送信する。通信インタフェース506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506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。

40

【0044】

アクセスノード600、管理ノード30および端末装置21,22,23も、コントロールノード500と同様のハードウェア構成によって実現できる。ただし、アクセスノード600は、ネットワーク10と接続するための通信インタフェースに加えて、ネットワーク20と接続するための通信インタフェースを更に備えている。

【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。

ここで、コントロールノード500がアクセスノード600に対して提供する論理ボリュームについて説明する。論理ボリュームは、ストレージノード100,200,300

50

、400によって分散管理されているデータを、アクセスノード600から容易に利用できるようにするための仮想的なボリュームである。

【0046】

図5は、論理ボリュームのデータ構造を示す第1の模式図である。論理ボリューム700には、“VV-A”という論理ボリュームIDが付与されている。また、ストレージノード100、200、300、400には、それぞれ“SN-A”、“SN-B”、“SN-C”、“SN-D”というノードIDが付与されている。

【0047】

各ストレージノード100、200、300、400に接続されたストレージ装置110、210、310、410それぞれにおいてRAID5の論理ディスクが構成されている。この論理ディスクは6つのスライスに分割されて、個々のストレージノード内で管理されている。

10

【0048】

図5の例では、ストレージ装置110内の記憶領域は、6つのスライス121～126に分割されている。ストレージ装置210内の記憶領域は、6つのスライス221～226に分割されている。ストレージ装置310内の記憶領域は、6つのスライス321～326に分割されている。ストレージ装置410内の記憶領域は、6つのスライス421～426に分割されている。

【0049】

論理ボリューム700は、セグメント710、720、730、740、750、760という単位で構成される。セグメント710、720、730、740、750、760は、それぞれプライミスライス711、721、731、741、751、761とセカンダリスライス712、722、732、742、752、762との組で構成される。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ように配置を行う。

20

【0050】

図5の例では、スライスIDを、“P”または“S”のアルファベットと数字との組み合わせで示している。“P”はプライミ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710のプライミスライスが“P1”で示され、セカンダリスライスが“S1”で示される。

30

【0051】

このような構造の論理ボリューム700の各プライミスライスおよびセカンダリスライスが、ストレージ装置110、210、310、410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710のプライミスライス711は、ストレージ装置110のスライス121に対応付けられ、セカンダリスライス712は、ストレージ装置310のスライス322に対応付けられている。

【0052】

そして、ストレージ装置110、210、310、410には、自己のスライスに対応付けられたプライミスライスまたはセカンダリスライスのデータが格納される。

40

次に、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のモジュール構成について説明する。

【0053】

図6は、ストレージノードの機能を示すブロック図である。ストレージノード100は、スライス情報記憶部130、データアクセス部140、スライス情報管理部150および負荷監視部160を有する。

【0054】

スライス情報記憶部130には、ストレージ装置110が有するスライスについてのスライス情報が格納される。スライス情報には、スライスを特定するためのアドレスやスライスへの割り当ての種別（プライミスライスまたはセカンダリスライスのいずれか）な

50

どの情報が含まれている。また、プライマリスライスについては、対応するセカンダリスライスを管理しているストレージノードに関する情報も含まれている。

【 0 0 5 5 】

データアクセス部 1 4 0 は、アクセスノード 6 0 0 からのアクセスを受け付けると、スライス情報記憶部 1 3 0 に格納されたスライス情報を参照して、ストレージ装置 1 1 0 に格納されたデータを操作する。

【 0 0 5 6 】

具体的には、データアクセス部 1 4 0 は、アクセスノード 6 0 0 からアドレスを指定したデータの読み込み要求 (R e a d 要求) を受け付けると、指定されたアドレスが属するスライスがプライマリスライスであるか否か判断する。プライマリスライスである場合には、データアクセス部 1 4 0 は、指定されたアドレスに対応するデータをストレージ装置 1 1 0 から取得し、アクセスノード 6 0 0 に送信する。プライマリスライスでない場合には、データアクセス部 1 4 0 は、アドレスの指定が不適切である旨をアクセスノード 6 0 0 に通知する。

10

【 0 0 5 7 】

また、データアクセス部 1 4 0 は、アクセスノード 6 0 0 からアドレスおよび書き込み内容を指定したデータの書き込み要求 (W r i t e 要求) を受け付けると、ストレージ装置 1 1 0 内の指定されたアドレスに対応する位置にデータの書き込みを試みる。そして、データアクセス部 1 4 0 は、書き込みの結果を、アクセスノード 6 0 0 に通知する。ここで、指定されたアドレスが属するスライスがプライマリスライスである場合には、対応するセカンダリスライスを管理するストレージノードに対して、同様の書き込みを行うよう指示する。これによって、プライマリスライスとセカンダリスライスとの内容が一致するように維持される。

20

【 0 0 5 8 】

スライス情報管理部 1 5 0 は、ストレージノード 1 0 0 の稼働状態を定期的にコントロールノード 5 0 0 に通知する。また、スライス情報管理部 1 5 0 は、コントロールノード 5 0 0 からスライス情報の取得要求があると、スライス情報記憶部 1 3 0 に格納されたスライス情報を送信する。また、スライス情報管理部 1 5 0 は、スライス情報の更新の指示を受け付けると、指示された更新内容を、スライス情報記憶部 1 3 0 に格納されたスライス情報に反映させる。

30

【 0 0 5 9 】

また、スライス情報管理部 1 5 0 は、ストレージ装置 1 1 0 が H D D へのアクセスを効率化するためにデータ格納位置の再編成を開始すると、その旨をコントロールノード 5 0 0 に通知する。再編成が終了すると、再編成後の各スライスのアドレスを、スライス情報記憶部 1 3 0 に格納されたスライス情報に反映させると共に、更新後のスライス情報をコントロールノード 5 0 0 に対して送信する。

【 0 0 6 0 】

負荷監視部 1 6 0 は、ストレージノード 1 0 0 の負荷を継続的に監視する。例えば、負荷監視部 1 6 0 は、C P U 1 0 1 の稼働率、H D D インタフェース 1 0 3 に対する I / O (Input / Output) 要求の回数、通信インタフェース 1 0 6 の使用率などを監視する。そして、負荷監視部 1 6 0 は、ストレージノード 1 0 0 の負荷を示す負荷情報を、定期的にコントロールノード 5 0 0 に対して送信する。なお、スライス情報管理部 1 5 0 による稼働状態の通知と負荷監視部 1 6 0 による負荷情報の送信とを兼ねるように、コントロールノード 5 0 0 に対して通信を行うことも可能である。

40

【 0 0 6 1 】

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。

図 7 は、コントロールノードおよびアクセスノードの機能を示すブロック図である。

【 0 0 6 2 】

コントロールノード 5 0 0 は、論理ボリューム記憶部 5 1 0 、負荷情報記憶部 5 2 0 、

50

論理ボリューム管理部 530、負荷情報収集部 540 および種別変更部 550 を有する。

論理ボリューム記憶部 510 には、1つ以上の論理ボリュームが格納される。論理ボリュームでは、ストレージ装置 110, 210, 310, 410 が管理する記憶領域を一元的に扱うために、仮想的なアドレスである論理アドレスによって各セグメントが管理されている。論理ボリュームには、セグメントを特定するための論理アドレスやセグメントに属するプライマリスライスおよびセカンダリスライスを特定するための情報などが含まれている。

【0063】

負荷情報記憶部 520 には、ストレージノード 100, 200, 300, 400 それぞれの最新の負荷を示す負荷情報が格納される。負荷情報記憶部 520 に格納される負荷情報は、ストレージノード 100, 200, 300, 400 から取得したものである。

10

【0064】

論理ボリューム管理部 530 は、ストレージノード 100, 200, 300, 400 から稼働状態を示す通知をネットワーク 10 経由で受信する。これにより、論理ボリューム管理部 530 は、ストレージノード 100, 200, 300, 400 が正常に稼働しているか否かを知る。また、論理ボリューム管理部 530 は、必要に応じてストレージノード 100, 200, 300, 400 からスライス情報を取得し、論理ボリューム記憶部 510 に格納された論理ボリュームを更新する。また、論理ボリューム管理部 530 は、論理ボリューム記憶部 510 が更新されると、更新によって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。

20

【0065】

また、論理ボリューム管理部 530 は、ストレージノード 100, 200, 300, 400 から再編成の開始の通知を受け付けると、負荷情報記憶部 520 に再編成中であることを示す情報を格納する。これは、再編成が開始されると、そのストレージ装置の負荷が高くなり、アクセスに対する応答性能が極端に低下することを考慮したものである。そして、論理ボリューム管理部 530 は、再編成後のスライス情報をストレージノード 100, 200, 300, 400 から受信すると、受信したスライス情報を論理ボリューム記憶部 510 に反映させると共に、負荷情報記憶部 520 から再編成中であることを示す情報を削除する。

【0066】

30

また、論理ボリューム管理部 530 は、アクセスノード 600 から論理ボリュームの参照要求を受け付けると、論理ボリューム記憶部 510 に格納された論理ボリュームを、アクセスノード 600 に対して送信する。

【0067】

負荷情報収集部 540 は、ストレージノード 100, 200, 300, 400 が定期的には送信する負荷情報を、ネットワーク 10 経由で受信する。そして、負荷情報収集部 540 は、受信した負荷情報を負荷情報記憶部 520 に格納する。

【0068】

種別変更部 550 は、論理ボリューム記憶部 510 に格納された論理ボリュームと負荷情報記憶部 520 に格納された負荷情報とを継続的に監視し、プライマリスライスとセカンダリスライスとを交換すべきセグメントを判定する。具体的には、種別変更部 550 は、プライマリスライスを割り当てたストレージノードの負荷とセカンダリスライスを割り当てたストレージノードの負荷との差が所定の許容量を超える場合に、そのプライマリスライスとセカンダリスライスとを交換すべきと判定する。そして、種別変更部 550 は、判定結果に応じて論理ボリューム記憶部 510 に格納された論理ボリュームを更新する。判定方法の詳細については、後で詳細に説明する。

40

【0069】

なお、論理ボリューム管理部 530 および負荷情報収集部 540 は、ストレージノード 100, 200, 300, 400 から各種情報が送られてくるのを待つ代わりに、定期的にストレージノード 100, 200, 300, 400 にアクセスして各種情報を収集する

50

ようにしてもよい。

【0070】

アクセスノード600は、論理ボリューム記憶部610およびデータアクセス制御部620を有する。

論理ボリューム記憶部610には、コントロールノード500の論理ボリューム記憶部510に格納されている論理ボリュームと同じ情報が格納される。

【0071】

データアクセス制御部620は、実行中のプログラムからデータへのアクセス要求があると、まず論理ボリューム記憶部610に論理ボリュームが格納されているか否か確認する。論理ボリュームが格納されていない場合、データアクセス制御部620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を論理ボリューム記憶部610に格納する。

10

【0072】

そして、データアクセス制御部620は、論理ボリュームに基づいてアクセス先のストレージノードを特定する。すなわち、利用するデータが属するセグメントを特定し、特定したセグメントのプライマリスライスを管理するストレージノードを特定する。その後、データアクセス制御部620は、特定したストレージノードにアクセスする。ここで、アクセスに失敗した場合、コントロールノード500から論理ボリュームを取得した後にデータの配置状況が変化していることが考えられるため、データアクセス制御部620は、コントロールノード500から最新の論理ボリュームを取得して、再びストレージノードへのアクセスを試みる。

20

【0073】

図8は、スライス情報テーブルのデータ構造例を示す図である。図8に示すスライス情報テーブル131は、ストレージノード100のスライス情報記憶部130に格納されている。スライス情報テーブル131には、ディスクを示す項目、物理アドレスを示す項目、ブロック数を示す項目、ボリュームを示す項目、論理アドレスを示す項目、種別を示す項目およびセカンダリを示す項目が設けられている。各項目の横方向に並べられた情報同士が互いに関連付けられて、1つのスライスについてのスライス情報を構成する。

【0074】

ディスクを示す項目には、HDDを識別するディスクIDが設定される。物理アドレスを示す項目には、スライスの先頭ブロックを示す物理的なアドレスが設定される。ブロック数を示す項目には、スライスに含まれるブロックの数が設定される。

30

【0075】

ボリュームを示す項目には、スライスに対応付けられたセグメントが属する論理ボリュームの論理ボリュームIDが設定される。論理アドレスを示す項目には、スライスに対応付けられたセグメントの先頭の論理アドレスが設定される。種別を示す項目には、“P”または“S”のいずれかの値が設定される。“P”はプライマリスライスを意味し、“S”はセカンダリスライスを意味する。

【0076】

セカンダリを示す項目には、種別がプライマリスライスである場合、対応するセカンダリスライスの割り当て先の情報が設定される。具体的には、ストレージノードのノードID、HDDのディスクID、スライスの先頭ブロックを示す物理アドレスが設定される。種別がセカンダリスライスである場合、セカンダリを示す項目は空欄となる。

40

【0077】

スライス情報テーブル131に格納されるスライス情報は、スライス情報管理部150によって適宜更新される。例えば、ディスクが“sd-a”、物理アドレスが“3072”、ブロック数が“512”、ボリュームが“VV-1”、論理アドレスが“4096”、種別が“S”という情報が格納される。これは、ディスクID“sd-a”のディスクの3072ブロックから3583ブロックまでの記憶領域が1つのスライスを構成し、そのスライスに、論理アドレスが4096ブロックから4607ブロックまでのセグメント

50

がセカンダリスライスとして割り当てられていることを示している。

【 0 0 7 8 】

図 9 は、ボリューム情報テーブルのデータ構造例を示す図である。図 9 に示す論理ボリュームテーブル 5 1 1 は、論理ボリューム ID が “ V V - 1 ” の論理ボリュームについてのテーブルである。論理ボリュームテーブル 5 1 1 は、コントロールノード 5 0 0 の論理ボリューム記憶部 5 1 0 に格納されている。論理ボリュームテーブル 5 1 1 には、セグメントを示す項目、論理アドレスを示す項目、ブロック数を示す項目、種別を示す項目、ノードを示す項目、ディスクを示す項目および物理アドレスを示す項目が設けられている。各項目の横方向に並べられた情報同士が互いに関連付けられている。

【 0 0 7 9 】

セグメントを示す項目には、セグメントを識別するセグメント ID が設定される。論理アドレスを示す項目には、セグメントの先頭の論理アドレスが設定される。ブロック数を示す項目には、セグメントに含まれるブロックの数が設定される。

【 0 0 8 0 】

種別を示す項目には、“ P ” または “ S ” のいずれかの値が設定される。ノードを示す項目には、割り当て先のストレージノードを識別するノード ID が設定される。ディスクを示す項目には、ストレージノード内で HDD を識別するディスク ID が設定される。物理アドレスを示す項目には、割り当て先のスライスの先頭ブロックを示す物理的なアドレスが設定される。

【 0 0 8 1 】

論理ボリュームテーブル 5 1 1 に格納される情報は、ストレージノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 から取得したスライス情報に基づいて、論理ボリューム管理部 5 3 0 によって生成される。

【 0 0 8 2 】

図 1 0 は、負荷情報テーブルのデータ構造例を示す図である。図 1 0 に示す負荷情報テーブル 5 2 1 は、コントロールノード 5 0 0 の負荷情報記憶部 5 2 0 に格納されている。負荷情報テーブル 5 2 1 には、ノードを示す項目、CPU 負荷を示す項目、I / O 数を示す項目、期間 I / O 数を示す項目および I / F 使用率を示す項目が設けられている。各項目の横方向に並べられた情報同士が互いに関連付けられて、1 つのストレージノードについての負荷情報を構成する。

【 0 0 8 3 】

ノードを示す項目には、ストレージノードを識別するノード ID が設定される。CPU 負荷を示す項目には、CPU の現在の稼働率が設定される。I / O 数を示す項目には、過去に発生した I / O の総回数が設定される。期間 I / O 数を示す項目には、所定期間内に発生した I / O の回数が設定される。例えば、直近 2 4 時間に発生した I / O の回数が設定される。I / F 使用率を示す項目には、ネットワークの現在の使用率が設定される。

【 0 0 8 4 】

負荷情報テーブル 5 2 1 に格納される負荷情報は、負荷情報収集部 5 4 0 によって適宜更新される。例えば、ノードが “ S N - A ”、CPU 負荷が “ 7 5 % ”、I / O 数が “ 1 0 0 0 0 ”、期間 I / O 数が “ 1 0 0 0 ”、I / F 使用率が “ 6 6 % ” という情報が格納される。

【 0 0 8 5 】

次に、以上のような構成およびデータ構造のシステムにおいて実行される処理の詳細を説明する。最初に、コントロールノード 5 0 0 の種別変更部 5 5 0 が、ストレージノードの負荷に応じてプライマリスライスとセカンダリスライスとを交換する処理について説明する。

【 0 0 8 6 】

図 1 1 は、第 1 の種別変更処理の手順を示すフローチャートである。この種別変更処理は、種別変更部 5 5 0 によって定期的に行われる。以下、図 1 1 に示す処理をステップ番号に沿って説明する。

10

20

30

40

50

【0087】

[ステップS11] 種別変更部550は、負荷情報記憶部520に格納された負荷情報テーブル521を参照して、負荷が最も高いストレージノードを特定する。

[ステップS12] 種別変更部550は、負荷情報テーブル521を参照して、負荷が最も低いストレージノードを特定する。

【0088】

[ステップS13] 種別変更部550は、ステップS11で特定した最高負荷のストレージノードの負荷と、ステップS12で特定した最低負荷のストレージノードの負荷との差が、所定の閾値を超えるか否か判断する。閾値を超える場合には、処理がステップS14に進められる。閾値以下である場合には、処理が終了する。

10

【0089】

[ステップS14] 種別変更部550は、論理ボリューム記憶部510に格納された論理ボリュームテーブル511を参照して、ステップS11で特定した最高負荷のストレージノードにプライマリスライスが割り当てられ、ステップS12で特定した最低負荷のストレージノードのセカンダリスライスが割り当てられている、セグメントが存在するか否か判断する。存在する場合には、処理がステップS15に進められる。存在しない場合には、処理が終了する。

【0090】

[ステップS15] 種別変更部550は、論理ボリュームテーブル511を更新して、ステップS14の判断条件を具備するプライマリスライスとセカンダリスライスとの間で種別を交換する。

20

【0091】

[ステップS16] 論理ボリューム管理部530は、論理ボリュームテーブル511が更新されたことを検知し、ステップS15で交換されたプライマリスライスおよびセカンダリスライスそれぞれを管理する2つのストレージノードに対して更新内容を通知する。

【0092】

ここで、負荷を比較する際に、CPU負荷、I/O数、期間I/O数およびI/F利用率のいずれの指標を用いるかは、分散ストレージシステムの管理者によって予め種別変更部550に設定されている。管理者は、いずれか1つの指標を用いて比較を行うように設定してもよいし、2つ以上の指標から所定の計算式によって計算される値を用いて比較を行うように設定してもよい。

30

【0093】

なお、図11に示した方法では、プライマリスライスを最高負荷のストレージノードから選択し、セカンダリスライスを最低負荷のストレージノードから選択することとしたが、セカンダリスライスを、最低負荷のストレージノード以外のストレージノードからも選択できるようにすることも考えられる。すなわち、プライマリスライスを最高負荷のストレージノードから選択し、セカンダリスライスを最高負荷のストレージノードとの負荷の差が許容量を超えている全てのストレージノードから選択できるようにしてもよい。このとき、複数のセカンダリスライスの候補がある場合には、最高負荷のストレージノードとの負荷の差ができる限り大きいストレージノードから選択するようにすればよい。

40

【0094】

このようにして、種別変更部550は、最高負荷のストレージノードと最低負荷のストレージノードとの負荷の差が許容量を超えており、最高負荷のストレージノードにプライマリスライスが割り当てられ最低負荷のストレージノードにセカンダリスライスが割り当てられているセグメントが存在する場合には、両者の種別を交換する。これにより、最高負荷のストレージノードに対するアクセスの少なくとも一部が、最低負荷のストレージノードに振り向けられる。このとき、プライマリスライスとセカンダリスライスの内容は同一であるため、論理ボリュームを更新するのみでよく、ストレージ装置110, 210, 310, 410に格納されたデータ自体を移動する必要はない。

【0095】

50

交換対象となるプライミスライスとセカンダリスライスとを検出する方法は、上記の方法以外にも、さまざまな方法が考えられる。次に、他の検出方法の例を示す。

図12は、第2の種別変更処理の手順を示すフローチャートである。この種別変更処理は、種別変更部550によって定期的に行われる。以下、図12に示す処理をステップ番号に沿って説明する。

【0096】

[ステップS21] 種別変更部550は、論理ボリューム記憶部510に格納された論理ボリュームテーブル511を参照して、ストレージノードを1つ選択する。

[ステップS22] 種別変更部550は、論理ボリュームを参照して、ステップS21で選択したストレージノードにプライミスライスまたはセカンダリスライスを割り当てているセグメントを1つ選択する。

10

【0097】

[ステップS23] 種別変更部550は、負荷情報記憶部520に格納された負荷情報テーブル521を参照して、ステップS22で選択したセグメントのプライミスライスの割り当て先のストレージノードとセカンダリスライスの割り当て先のストレージノードとの間で負荷の差を計算する。なお、セカンダリスライスの割り当て先のストレージノードの方が負荷が高い場合には、負荷の差は0とする。

【0098】

[ステップS24] 種別変更部550は、ステップS22で条件を具備する全てのセグメントを選択したか否か判断する。全てのセグメントを選択した場合には、処理がステップS25に進められる。未選択のセグメントが存在する場合には、処理がステップS22に進められる。

20

【0099】

[ステップS25] 種別変更部550は、ステップS23で計算した負荷の差の最大値が、所定の閾値を超えるか否か判断する。閾値を超える場合には、処理がステップS26に進められる。閾値以下である場合には、処理がステップS28に進められる。

【0100】

[ステップS26] 種別変更部550は、論理ボリュームテーブル511を更新して、ステップS23で計算した負荷の差が最大となったセグメントのプライミスライスとセカンダリスライスとの間で種別を交換する。

30

【0101】

[ステップS27] 論理ボリューム管理部530は、論理ボリュームテーブル511が更新されたことを検知し、ステップS26で交換されたプライミスライスおよびセカンダリスライスそれぞれを管理する2つのストレージノードに対して更新内容を通知する。

【0102】

[ステップS28] 種別変更部550は、ステップS21で全てのストレージノードを選択したか否か判断する。全てのストレージノードを選択した場合には、処理が終了する。未選択のストレージノードが存在する場合には、処理がステップS29に進められる。

【0103】

[ステップS29] 種別変更部550は、一定時間待機する。その後、処理がステップS21に進められる。

40

なお、上記ステップS21では、ストレージノードを選択する順序を、負荷が高い順または負荷が低い順とすることも考えられる。

【0104】

このようにして、種別変更部550は、ストレージノードを順次確認し、プライミスライスが割り当てられているストレージノードとセカンダリスライスが割り当てられているストレージノードとの負荷の差が許容量を超えているセグメントが存在する場合には、両者の種別を交換する。これにより、最高負荷のストレージノードにプライミスライスが割り当てられ最低負荷のストレージノードにセカンダリスライスが割り当てられているセグメントが存在しない場合でも、負荷を分散させることができる。

50

【 0 1 0 5 】

図 1 3 は、論理ボリュームのデータ構造を示す第 2 の模式図である。図 1 3 に示すセグメントの割り当て状況は、上記の種別変更処理が実行されることで、図 5 に示したセグメントの割り当て状況が変化したものである。図 1 3 の例では、セグメント 7 5 0 のプライマリスライスとセカンダリスライスとが交換されている。すなわち、交換前は、プライマリスライス 7 5 1 がストレージ装置 1 1 0 に割り当てられ、セカンダリスライス 7 5 2 がストレージ装置 4 1 0 に割り当てられていたが、交換後は、プライマリスライス 7 5 1 がストレージ装置 4 1 0 に割り当てられ、セカンダリスライス 7 5 2 がストレージ装置 1 1 0 に割り当てられている。これにより、ストレージノード 1 0 0 の負荷が軽減される。

【 0 1 0 6 】

次に、再編成中のストレージ装置にプライマリスライスを割り当てないようにして、応答性能の低下を避ける処理について説明する。

図 1 4 は、種別一斉変更処理の手順を示すフローチャートである。以下、図 1 4 に示す処理をステップ番号に沿って説明する。

【 0 1 0 7 】

[ステップ S 3 1] 論理ボリューム管理部 5 3 0 は、ストレージノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 のいずれかから再編成の開始を示す通知を受け付けると、負荷情報記憶部 5 2 0 に再編成の開始を示す情報を格納する。

【 0 1 0 8 】

[ステップ S 3 2] 種別変更部 5 5 0 は、論理ボリューム記憶部 5 1 0 に格納された論理ボリュームテーブル 5 1 1 を参照して、再編成中のストレージノードにプライマリスライスを割り当てているセグメントを 1 つ選択する。

【 0 1 0 9 】

[ステップ S 3 3] 種別変更部 5 5 0 は、論理ボリュームテーブル 5 1 1 を更新して、ステップ S 3 2 で選択したセグメントのプライマリスライスとセカンダリスライスとの間で種別を交換する。

【 0 1 1 0 】

[ステップ S 3 4] 論理ボリューム管理部 5 3 0 は、論理ボリュームテーブル 5 1 1 が更新されたことを検知し、ステップ S 3 3 で交換されたプライマリスライスおよびセカンダリスライスそれぞれを管理する 2 つのストレージノードに対して更新内容を通知する。

【 0 1 1 1 】

[ステップ S 3 5] 種別変更部 5 5 0 は、ステップ S 3 2 で条件を具備する全てのセグメントを選択したか否か判断する。全てのセグメントを選択した場合には、処理が終了する。未選択のセグメントが存在する場合には、処理がステップ S 3 2 に進められる。

【 0 1 1 2 】

このようにして、種別変更部 5 5 0 は、再編成中のストレージノードに割り当てられたプライマリスライスを全てセカンダリスライスに変更する。これにより、再編成を行うことによる応答性能の低下を避けることができる。また、再編成が終了すると、その後は負荷状態に応じてプライマリスライスとセカンダリスライスとの交換が行われ、再編成が終了したストレージノードに再びアクセスの一部が振り向けられる。

【 0 1 1 3 】

次に、アクセスノード 6 0 0 がストレージノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 にアクセスする際の通信の流れの例について説明する。

図 1 5 は、アクセスノードからのデータアクセスの流れを示すシーケンス図である。以下、図 1 5 に示す処理をステップ番号に沿って説明する。

【 0 1 1 4 】

[ステップ S 4 1] アクセスノード 6 0 0 は、コントロールノード 5 0 0 に参照要求を送信する。

[ステップ S 4 2] コントロールノード 5 0 0 は、参照要求に応答して、最新の論理ボリュームをアクセスノード 6 0 0 に送信する。

10

20

30

40

50

【 0 1 1 5 】

【ステップS 4 3】アクセスノード6 0 0は、ステップS 4 2で取得した論理ボリュームに基づいてストレージノードにアクセスする。ここでは、アクセスノード6 0 0は、ストレージノード1 0 0にR e a d要求を送信する。

【 0 1 1 6 】

【ステップS 4 4】ストレージノード1 0 0は、R e a d要求に応答して、ストレージ装置1 1 0からデータを取得し、取得したデータをアクセスノード6 0 0に送信する。

【ステップS 4 5】コントロールノード5 0 0は、ストレージノードの負荷を分散させるために、一部のセグメントのプライマリスライスとセカンダリスライスとを交換する。ここでは、コントロールノード5 0 0は、ステップS 4 3でアクセスノード6 0 0が読み込みを行ったデータが属するセグメントのプライマリスライスとセカンダリスライスとを交換するものとする。

10

【 0 1 1 7 】

【ステップS 4 6 a】コントロールノード5 0 0は、交換前にプライマリスライスが割り当てられていたストレージノード1 0 0に対して、セカンダリスライスに種別が変更された旨を通知する。

【 0 1 1 8 】

【ステップS 4 6 b】コントロールノード5 0 0は、交換前にセカンダリスライスが割り当てられていたストレージノードに対して、プライマリスライスに種別が変更された旨を通知する。ここでは、交換前にセカンダリスライスが割り当てられていたストレージノードは、ストレージノード4 0 0であるとする。

20

【 0 1 1 9 】

【ステップS 4 7】アクセスノード6 0 0は、ステップS 4 2で取得した論理ボリュームに基づいて、ステップS 4 3で読み込みを行ったデータを再度読み込む。すなわち、アクセスノード6 0 0は、ストレージノード1 0 0にR e a d要求を送信する。

【 0 1 2 0 】

【ステップS 4 8】ストレージノード1 0 0は、R e a d要求で示されるデータが属するスライスがセカンダリスライスであるため、要求を拒否する旨をアクセスノード6 0 0に応答する。

【 0 1 2 1 】

【ステップS 4 9】アクセスノード6 0 0は、コントロールノード5 0 0に参照要求を送信する。

30

【ステップS 5 0】コントロールノード5 0 0は、参照要求に応答して、最新の論理ボリュームをアクセスノード6 0 0に送信する。

【 0 1 2 2 】

【ステップS 5 1】アクセスノード6 0 0は、ステップS 5 0で取得した論理ボリュームに基づいて、ストレージノード4 0 0にR e a d要求を送信する。

【ステップS 5 2】ストレージノード4 0 0は、R e a d要求に応答して、ストレージ装置4 1 0からデータを取得し、取得したデータをアクセスノード6 0 0に送信する。

【 0 1 2 3 】

このようにして、アクセスノード6 0 0は、コントロールノード5 0 0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先のストレージノードを特定する。ここで、コントロールノード5 0 0がプライマリスライスとセカンダリスライスとの交換を行うと、アクセスノード6 0 0はストレージノードへのアクセスに失敗する場合がある。この場合、アクセスノード6 0 0は、コントロールノード5 0 0から最新の論理ボリュームを取得し直し、ストレージノードへのアクセスを再度試みる。

40

【 0 1 2 4 】

このような分散ストレージシステムによれば、アクセス用のプライマリスライスが割り当てられたストレージノードとバックアップ用のセカンダリスライスが割り当てられたストレージノードとの間で負荷が許容値を超えると、プライマリスライスとセカンダリスラ

50

イスとが自動的に交換される。

【 0 1 2 5 】

これにより、データを移動することなく、当該データのアクセス先のストレージノードが変更される。このとき、アクセス単位でアクセス先が判定されるわけではない。従って、負荷分散の処理自体が原因で応答が遅くなることを防止し、分散ストレージシステムの応答性能を確実に高めることができる。また、データの配置管理は、ファイル単位で行われるのではなく、記憶領域の大きさに応じて分割したスライス単位で行われるため、配置管理の機構が簡潔になる。

【 0 1 2 6 】

また、再編成中のストレージ装置に対しては、プライマリスライスを割り当てないようにすることで、アクセスができる限り発生しないようにすることができる。これにより、応答性能の低下を防止できる。

【 0 1 2 7 】

以上、本発明のストレージ管理プログラム、ストレージ管理装置およびストレージ管理方法を図示の実施の形態に基づいて説明したが、本発明はこれに限定されるものではなく、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工程が付加されていてもよい。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。

【 0 1 2 8 】

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、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。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としては、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 2 9 】

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の記憶装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピュータにそのプログラムを転送することもできる。

【 0 1 3 0 】

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

【 0 1 3 1 】

上記については単に本発明の原理を示すものである。さらに、多数の変形、変更が当業者にとって可能であり、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、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る。

【 符号の説明 】

【 0 1 3 2 】

10

20

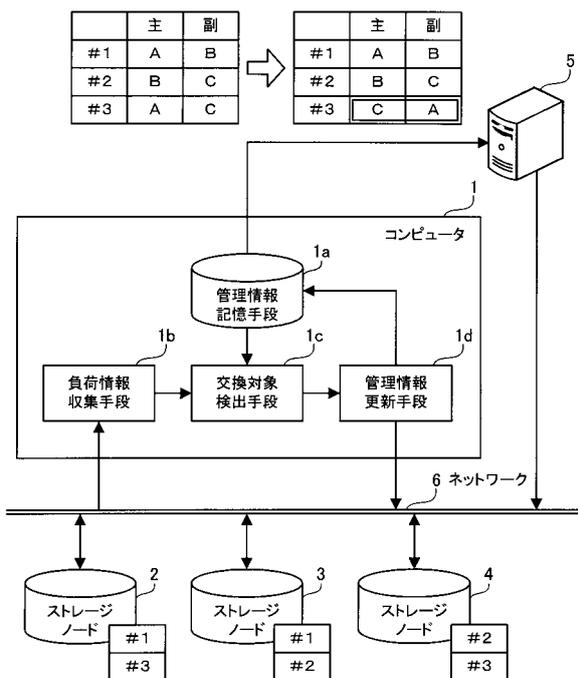
30

40

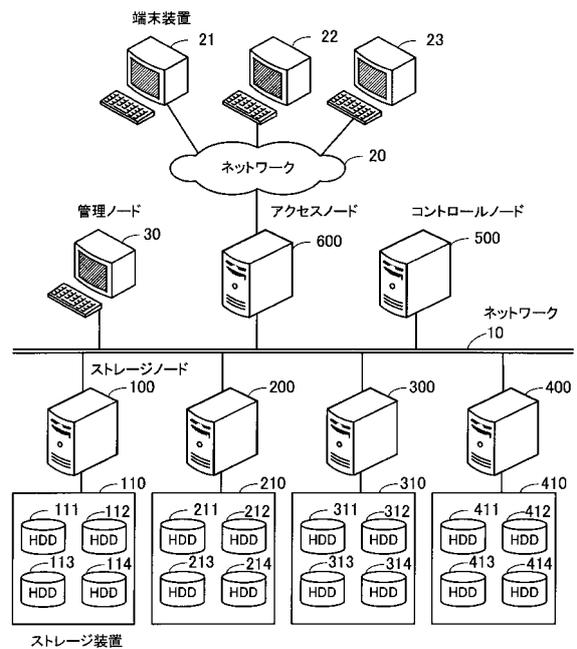
50

- 1, 5 コンピュータ
- 1 a 管理情報記憶手段
- 1 b 負荷情報収集手段
- 1 c 交換対象検出手段
- 1 d 管理情報更新手段
- 2, 3, 4 ストレージノード
- 6 ネットワーク

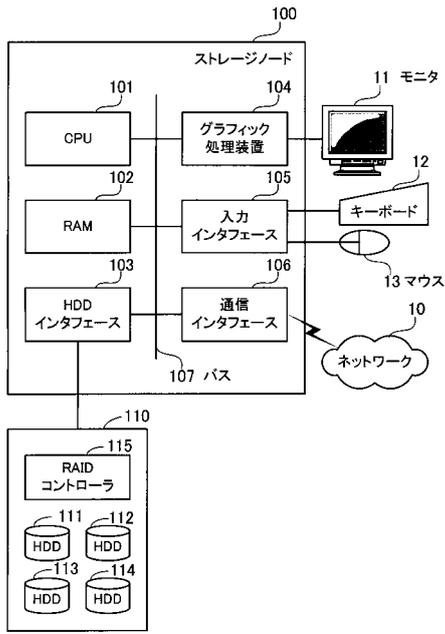
【図1】



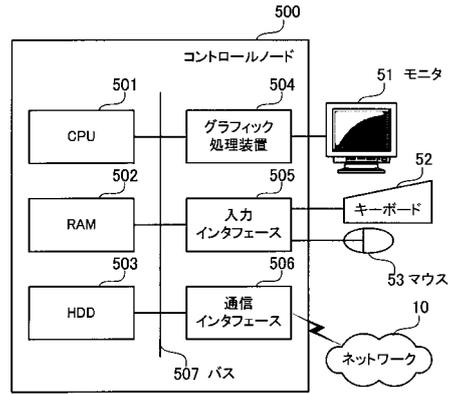
【図2】



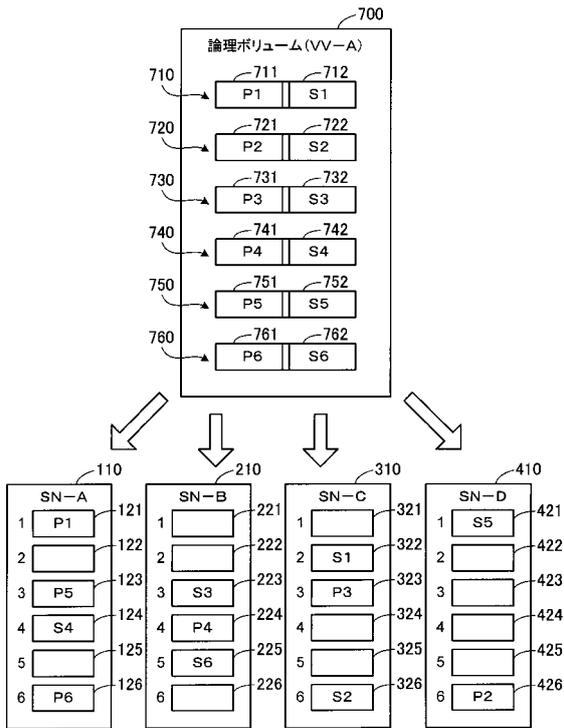
【図3】



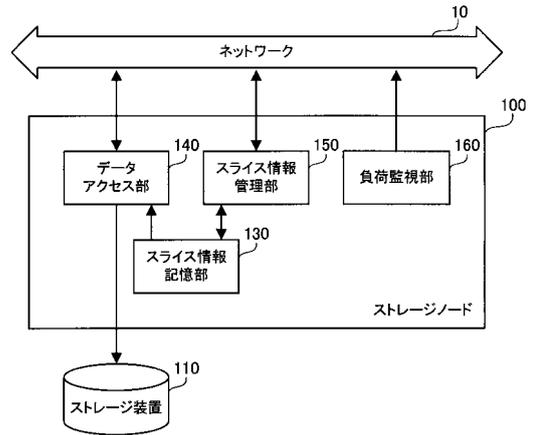
【図4】



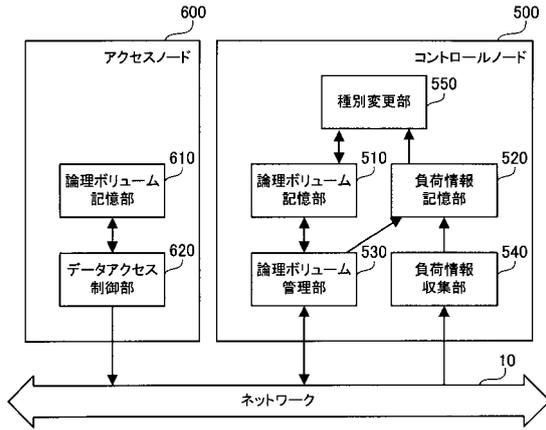
【図5】



【図6】



【図7】



【図8】

131

スライス情報テーブル						
ディスク	物理アドレス	ブロック数	ボリューム	論理アドレス	種別	セカンダリ
sd-a	0	1024	VV-1	0	P	SN-C, sd-a, 1024
	2048	512	VV-1	4608	P	SN-D, sd-b, 0
	3072	512	VV-1	4096	S	-
	5120	1024	VV-1	5120	P	SN-B, sd-b, 4096

【図9】

511

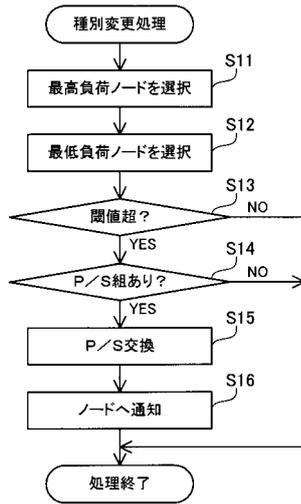
論理ボリュームテーブル (VV-1)						
セグメント	論理アドレス	ブロック数	種別	ノード	ディスク	物理アドレス
1	0	1024	P	SN-A	sd-a	0
			S	SN-C	sd-a	1024
2	1024	2048	P	SN-D	sd-b	5120
			S	SN-C	sd-a	5120
3	3072	1024	P	SN-C	sd-a	2048
			S	SN-B	sd-b	2048
4	4096	512	P	SN-B	sd-b	3072
			S	SN-A	sd-a	3072
5	4608	512	P	SN-A	sd-a	2048
			S	SN-D	sd-b	0
6	5120	1024	P	SN-A	sd-a	5120
			S	SN-B	sd-b	4096

【図10】

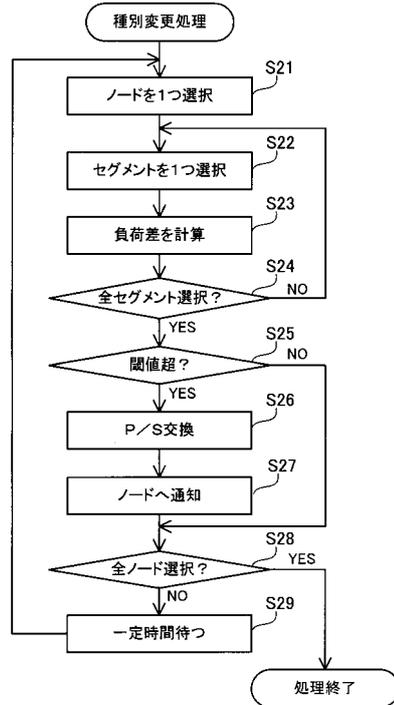
521

負荷情報テーブル				
ノード	CPU負荷	I/O数	期間 I/O数	I/F使用率
SN-A	75%	10000	1000	66%
SN-B	50%	5600	897	35%
SN-C	45%	6500	1210	10%
SN-D	30%	7880	3084	40%

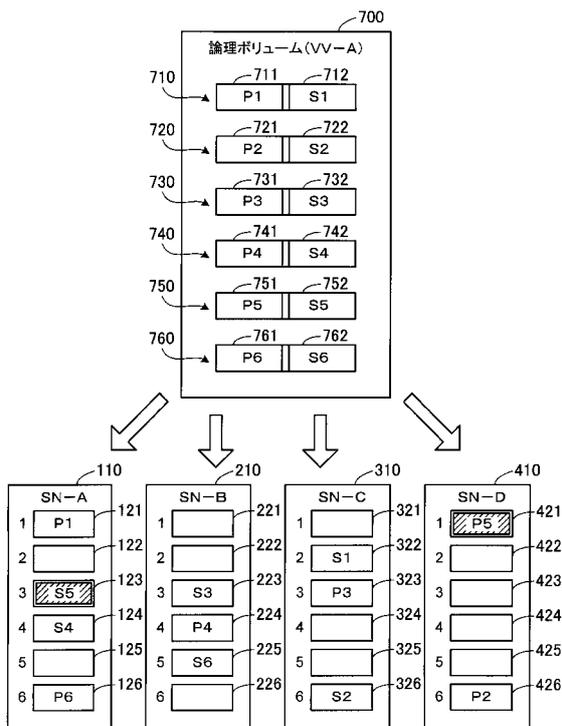
【図11】



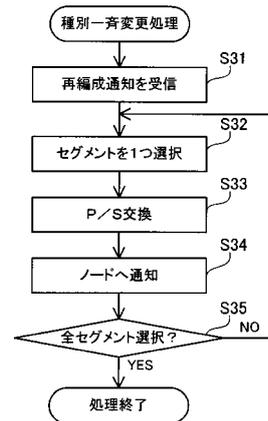
【図12】



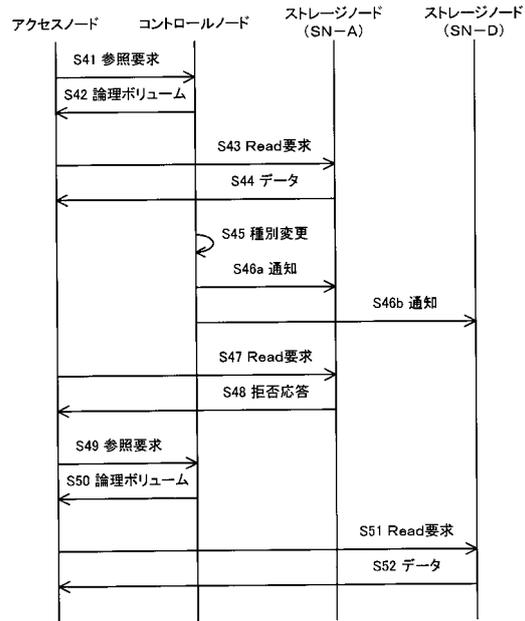
【図13】



【図14】



【図15】



フロントページの続き

- (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
- (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 大江 和一
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 渡辺 高志
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 熊野 達夫
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 北村 学

- (56)参考文献 特開2006-012005(JP,A)
特開2006-228188(JP,A)
特開2004-199264(JP,A)
特開平06-332782(JP,A)
特開平11-161555(JP,A)
特開2001-051963(JP,A)
特開2006-185413(JP,A)
特開2003-296167(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06

G06F 12/00

JSTPlus(JDreamII)