

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-20441
(P2010-20441A)

(43) 公開日 平成22年1月28日(2010.1.28)

(51) Int.Cl.
G06F 12/00 (2006.01)

F I
G06F 12/00 501B

テーマコード(参考)
5B082

審査請求 未請求 請求項の数 14 O L (全 37 頁)

(21) 出願番号 特願2008-178782(P2008-178782)
(22) 出願日 平成20年7月9日(2008.7.9)

(71) 出願人 000005108
株式会社日立製作所
東京都千代田区丸の内一丁目6番6号
(74) 代理人 100075513
弁理士 後藤 政喜
(74) 代理人 100114236
弁理士 藤井 正弘
(74) 代理人 100120260
弁理士 飯田 雅昭
(72) 発明者 田口 雄一
神奈川県川崎市麻生区王禅寺1099番地
株式会社日立製作所システム開発研究所
内

最終頁に続く

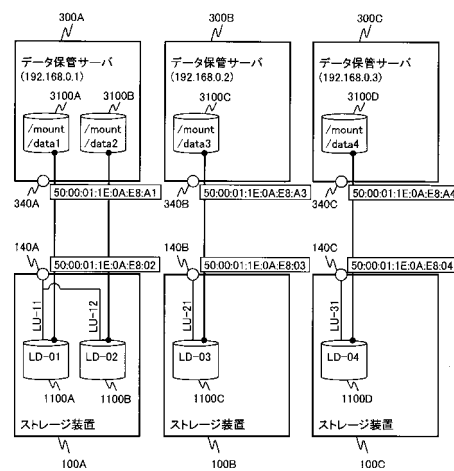
(54) 【発明の名称】 計算機システム、構成管理方法及び管理計算機

(57) 【要約】

【課題】 計算機の性能とストレージ装置の負荷とのバランスを改善するように構成を管理する。

【解決手段】 データ保管サーバと、データ保管サーバに記憶領域を提供するストレージ装置と、管理計算機と、を含む計算機システムであって、データ保管サーバは、記憶領域が割り当てられたデータ保管領域をクライアントに提供し、データと当該データを格納するデータ保管サーバとの対応を含む配置管理情報及びデータ保管領域と記憶領域との対応を含むデータ保管記憶領域構成情報を管理し、データの読み出し要求を受け付けた場合には、データ配置管理情報に基づいて要求されたデータが格納されたデータ保管サーバを特定し、特定されたデータ保管サーバから要求されたデータを取得し、管理計算機は、データ保管サーバ性能管理情報及び記憶領域稼働情報に基づいて、データ保管領域と記憶領域との対応を変更する。

【選択図】 図8



【特許請求の範囲】**【請求項 1】**

クライアントによってデータが格納される複数のデータ保管サーバと、前記データ保管サーバに記憶領域を提供するストレージ装置と、前記データ保管サーバ及び前記ストレージ装置にアクセス可能な管理計算機と、を含む計算機システムであって、

前記各データ保管サーバは、前記ストレージ装置に接続される第 1 インタフェースと、前記第 1 インタフェースに接続される第 1 プロセッサと、前記第 1 プロセッサに接続される第 1 メモリと、を備え、前記記憶領域が割り当てられたデータ保管領域を前記クライアントに提供し、

前記ストレージ装置は、前記データ保管サーバに接続される第 2 インタフェースと、前記第 2 インタフェースに接続される第 2 プロセッサと、前記第 2 プロセッサに接続される第 2 メモリと、を備え、

前記管理計算機は、前記データ保管サーバ及び前記ストレージ装置に接続される第 3 インタフェースと、前記第 3 インタフェースに接続される第 3 プロセッサと、前記第 3 プロセッサに接続される第 3 メモリと、を備え、

前記各データ保管サーバは、前記データと前記データが格納されるデータ保管サーバとの対応を含むデータ配置管理情報、及び前記データ保管領域と前記記憶領域との対応を含むデータ保管記憶領域構成情報を管理し、

前記データ配置管理情報は、前記データに基づいて決定されたデータ保管サーバに管理され、

前記管理計算機は、前記データ保管サーバの性能を含むデータ保管サーバ性能管理情報、及び前記記憶領域の負荷状態を含む記憶領域稼働情報を管理し、

前記データ保管サーバは、

前記クライアントからデータの読み出し要求を受け付けた場合には、前記要求されたデータに基づいて、前記要求されたデータに対応するデータ配置管理情報を管理するデータ保管サーバを特定し、

前記特定されたデータ保管サーバによって管理されるデータ配置管理情報を取得し、

前記取得されたデータ配置管理情報に基づいて、前記要求されたデータを格納するデータ保管サーバを特定し、

前記要求されたデータを格納するデータ保管サーバから、前記要求されたデータを取得し、

前記管理計算機は、

前記データ保管サーバ性能管理情報及び前記記憶領域稼働情報に基づいて、前記記憶領域との対応を変更するデータ保管領域を選択し、

前記記憶領域稼働情報に基づいて、前記選択されたデータ保管領域に新たに割り当てられる記憶領域を選択し、

前記選択されたデータ保管領域を提供するデータ保管サーバに、前記選択された新たに割り当てられる記憶領域を通知し、

前記データ保管サーバは、

前記管理計算機から通知された記憶領域を前記データ保管領域に割り当て、

前記データ配置管理情報及び前記データ保管記憶領域構成情報を更新することを特徴とする計算機システム。

【請求項 2】

前記管理計算機は、

前記記憶領域稼働情報に基づいて、最も負荷の高い記憶領域を選択し、

前記データ保管サーバ性能管理情報に基づいて、最も性能の高いデータ保管サーバを選択し、

前記選択されたデータ保管サーバによって提供されるデータ保管領域に、前記選択された記憶領域を割り当てることを特徴とする請求項 1 に記載の計算機システム。

【請求項 3】

前記管理計算機は、
前記記憶領域稼働情報に基づいて、最も負荷の低い記憶領域を選択し、
前記データ保管サーバ性能管理情報に基づいて、最も性能の低いデータ保管サーバを選択し、
前記選択されたデータ保管サーバによって提供されるデータ保管領域に、前記選択された記憶領域を割り当てることを特徴とする請求項 1 に記載の計算機システム。

【請求項 4】

前記ストレージ装置は、前記記憶領域に対するアクセスが許可されるデータ保管サーバの識別情報を含む記憶領域構成情報を管理し、

前記管理計算機は、前記対応が変更された記憶領域を提供するストレージ装置に、前記対応が変更されたデータ保管領域を提供するデータ保管サーバの識別情報を通知し、

前記ストレージ装置は、前記通知された識別情報によって識別されるデータ保管サーバによる、前記対応が変更された記憶領域へのアクセスを拒否することを特徴とする請求項 1 に記載の計算機システム。

【請求項 5】

前記データ保管サーバは、

前記データ保管領域と前記記憶領域との対応の変更を前記管理計算機から指示された場合には、前記データ配置管理情報を複製し、

前記複製されたデータ配置管理情報を更新し、

前記複製されたデータ配置管理情報を前記データ配置管理情報に反映させることを特徴とする請求項 1 に記載の計算機システム。

【請求項 6】

クライアントによってデータが格納される複数のデータ保管サーバと、前記データ保管サーバに記憶領域を提供するストレージ装置と、前記データ保管サーバ及び前記ストレージ装置にアクセス可能な管理計算機と、を含む計算機システムにおける構成管理方法であって、

前記各データ保管サーバは、前記ストレージ装置に接続される第 1 インタフェースと、前記第 1 インタフェースに接続される第 1 プロセッサと、前記第 1 プロセッサに接続される第 1 メモリと、を備え、前記記憶領域が割り当てられたデータ保管領域を前記クライアントに提供し、

前記ストレージ装置は、前記データ保管サーバに接続される第 2 インタフェースと、前記第 2 インタフェースに接続される第 2 プロセッサと、前記第 2 プロセッサに接続される第 2 メモリと、を備え、

前記管理計算機は、前記データ保管サーバ及び前記ストレージ装置に接続される第 3 インタフェースと、前記第 3 インタフェースに接続される第 3 プロセッサと、前記第 3 プロセッサに接続される第 3 メモリと、を備え、

前記各データ保管サーバは、前記データと前記データが格納されるデータ保管サーバとの対応を含むデータ配置管理情報、及び前記データ保管領域と前記記憶領域との対応を含むデータ保管記憶領域構成情報を管理し、

前記データ配置管理情報は、前記データに基づいて決定されたデータ保管サーバに管理され、

前記管理計算機は、前記データ保管サーバの性能を含むデータ保管サーバ性能管理情報、及び前記記憶領域の負荷状態を含む記憶領域稼働情報を管理し、

前記データ保管サーバは、

前記クライアントからデータの読み出し要求を受け付けた場合には、前記要求されたデータに基づいて、前記要求されたデータに対応するデータ配置管理情報を管理するデータ保管サーバを特定し、

前記特定されたデータ保管サーバによって管理されるデータ配置管理情報を取得し、

前記取得されたデータ配置管理情報に基づいて、前記要求されたデータを格納するデータ保管サーバを特定し、

10

20

30

40

50

前記要求されたデータを格納するデータ保管サーバから、前記要求されたデータを取得し、

前記方法は、

前記第3プロセッサが、前記データ保管サーバ性能管理情報及び前記記憶領域稼働情報に基づいて、前記記憶領域との対応を変更するデータ保管領域を選択し、

前記第3プロセッサが、前記記憶領域稼働情報に基づいて、前記選択されたデータ保管領域に新たに割り当てられる記憶領域を選択し、

前記第3プロセッサが、前記選択されたデータ保管領域を提供するデータ保管サーバに、前記選択された新たに割り当てられる記憶領域を通知し、

前記第1プロセッサが、前記管理計算機から通知された記憶領域を前記データ保管領域に割り当て、

前記第1プロセッサが、前記データ配置管理情報及び前記データ保管記憶領域構成情報を更新することを特徴とする構成管理方法。

【請求項7】

前記方法は、

前記第3プロセッサが、前記記憶領域稼働情報に基づいて、最も負荷の高い記憶領域を選択し、

前記第3プロセッサが、前記データ保管サーバ性能管理情報に基づいて、最も性能の高いデータ保管サーバを選択し、

前記第3プロセッサが、前記選択されたデータ保管サーバによって提供されるデータ保管領域に、前記選択された記憶領域を割り当てることを特徴とする請求項6に記載の構成管理方法。

【請求項8】

前記方法は、

前記第3プロセッサが、前記記憶領域稼働情報に基づいて、最も負荷の低い記憶領域を選択し、

前記第3プロセッサが、前記データ保管サーバ性能管理情報に基づいて、最も性能の低いデータ保管サーバを選択し、

前記第3プロセッサが、前記選択されたデータ保管サーバによって提供されるデータ保管領域に、前記選択された記憶領域を割り当てることを特徴とする請求項6に記載の構成管理方法。

【請求項9】

前記ストレージ装置は、前記記憶領域に対するアクセスが許可されたデータ保管サーバの識別情報を含む記憶領域構成情報を管理し、

前記方法は、

前記第3プロセッサが、前記対応が変更された記憶領域を提供するストレージ装置に、前記対応が変更されたデータ保管領域を提供するデータ保管サーバの識別情報を通知し、

前記第2プロセッサが、前記通知された識別情報によって識別されるデータ保管サーバによる、前記対応が変更された記憶領域へのアクセスを拒否することを特徴とする請求項6に記載の構成管理方法。

【請求項10】

前記方法は、

前記第1プロセッサが、前記データ保管領域と前記記憶領域との対応の変更を前記管理計算機から指示された場合には、前記データ配置管理情報を複製し、

前記第1プロセッサが、前記複製されたデータ配置管理情報を更新し、

前記第1プロセッサが、前記複製されたデータ配置管理情報を前記データ配置管理情報に反映させることを特徴とする請求項6に記載の構成管理方法。

【請求項11】

クライアントによってデータが格納される複数のデータ保管サーバと、前記データ保管サーバに記憶領域を提供するストレージ装置と、を含む計算機システムにおいて、前記デ

10

20

30

40

50

ータ保管サーバ及び前記ストレージ装置にアクセス可能な管理計算機であって、

前記データ保管サーバ及び前記ストレージ装置に接続されるインタフェースと、前記インタフェースに接続されるプロセッサと、前記プロセッサに接続されるメモリと、を備え、

前記各データ保管サーバは、前記記憶領域が割り当てられたデータ保管領域を前記クライアントに提供し、

前記プロセッサは、

前記データ保管サーバの性能を含むデータ保管サーバ性能管理情報、及び前記記憶領域の負荷状態を含む記憶領域稼働情報を管理し、

前記データ保管サーバ性能管理情報及び前記記憶領域稼働情報に基づいて、前記記憶領域との対応を変更するデータ保管領域を選択し、

前記記憶領域稼働情報に基づいて、前記選択されたデータ保管領域に新たに割り当てられる記憶領域を選択し、

前記選択されたデータ保管領域を提供するデータ保管サーバに、前記選択された新たに割り当てられる記憶領域を通知することを特徴とする管理計算機。

【請求項 1 2】

前記プロセッサは、

前記記憶領域稼働情報に基づいて、最も負荷の高い記憶領域を選択し、

前記データ保管サーバ性能管理情報に基づいて、最も性能の高いデータ保管サーバを選択し、

前記選択されたデータ保管サーバによって提供されるデータ保管領域に、前記選択された記憶領域を割り当てることを特徴とする請求項 1 1 に記載の管理計算機。

【請求項 1 3】

前記プロセッサは、

前記記憶領域稼働情報に基づいて、最も負荷の低い記憶領域を選択し、

前記データ保管サーバ性能管理情報に基づいて、最も性能の低いデータ保管サーバを選択し、

前記選択されたデータ保管サーバによって提供されるデータ保管領域に、前記選択された記憶領域を割り当てることを特徴とする請求項 1 1 に記載の管理計算機。

【請求項 1 4】

前記ストレージ装置は、前記記憶領域に対するアクセスが許可されたデータ保管サーバの識別情報を含む記憶領域構成情報を管理し、

前記プロセッサは、

前記対応が変更された記憶領域を提供するストレージ装置に、前記対応が変更されたデータ保管領域を提供するデータ保管サーバの識別情報を通知し、

前記通知された識別情報によって識別されるデータ保管サーバによる、前記対応が変更された記憶領域へのアクセスを拒否するように指示することを特徴とする請求項 1 1 に記載の管理計算機。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、計算機及び記憶装置によって構成される計算機システムに関し、特に、計算機システムの構成を管理する技術に関する。

【背景技術】

【0002】

一台以上の外部記憶装置と一台以上の計算機とを接続するストレージエリアネットワーク (SAN: Storage Area Network) が知られている。ストレージエリアネットワークは、複数の計算機が一つの大規模記憶装置を共有する場合に特に有効である。ストレージエリアネットワークを含むストレージシステムは、記憶装置又は計算機を容易に追加又は削除できるため、拡張性に優れている。

10

20

30

40

50

【0003】

また、SANに接続する外部記憶装置には、ディスクアレイ装置が一般的によく利用される。ディスクアレイ装置は、ハードディスクに代表される記憶デバイス（例えば、磁気記憶デバイス）を多数搭載する装置である。

【0004】

ディスクアレイ装置は、RAID（Redundant Array of Independent Disks）技術によって、数台の磁気記憶デバイスを一つのRAIDグループとして管理する。RAIDグループは、一つ以上の論理的な記憶領域を形成する。SANに接続された計算機は、当該記憶領域に対してデータ入出力処理を実行する。ディスクアレイ装置は、当該記憶領域にデータを記録する場合に、RAIDグループを構成する磁気記憶デバイスに冗長データを記録する。当該冗長データによって、磁気記憶デバイスの一つが故障した場合であっても、データを復元することができる。

10

【0005】

また、近年、契約書及び公的文書などを長期的に保存するために、長期アーカイブ目的のデータ保管システムが普及している。このようなデータ保管システムは、データ保護（WORM）、データ配置管理、重複データ排除、及び高速検索などの機能を有し、長期に渡ってデータを適切に保存及び管理するための仕組みを提供する（特許文献1参照）。

【0006】

重複データ排除は、複数の計算機が同一のデータを重複して保持している場合、各計算機の参照先（記憶領域内のデータ格納位置）を同一にすることによって、保持するデータ量を削減することができる。

20

【0007】

データ保管システムでは、システム運用が長期にわたるため、後から計算機資源を追加したり、障害が発生した装置を新たな装置に交換したりすることが想定される。したがって、もともとシステムに含まれていた計算機と、新たに導入された計算機との間に性能差が発生する可能性がある。

【0008】

このとき、新たに導入された計算機が高性能であっても、もともとシステムに含まれている計算機と比較して管理するデータ量が少量となってしまう場合がある。また、アクセス頻度の多いデータが古い計算機に格納されていると、低性能の計算機に負荷が集中し、システム全体の性能が低下してしまう。また、重複排除などのデータ量削減処理に伴い、データごとのアクセス頻度が変動し、計算機間で性能の負荷均衡が崩れたり、低性能サーバに負荷が集中したりする可能性もある。

30

【0009】

そこで、複数の計算機を用いて負荷を分散し、性能を改善するために高負荷の計算機から低負荷の計算機にデータを移行する技術が開示されている（特許文献2及び特許文献3参照）。

【特許文献1】国際公開第2005/043323号パンフレット

【特許文献2】特開2006-228188号公報

【特許文献3】特開2004-139200号公報

40

【発明の開示】

【発明が解決しようとする課題】

【0010】

特許文献2及び特許文献3に開示された技術では、負荷を分散し、データ保管システムの性能を改善させるために、システムに格納されたデータをストレージ装置間で移行させる必要がある。しかし、大容量のデータを格納するデータ保管システムでは、データを移行させる場合には、データの移行処理自体が高負荷となるという課題がある。

【0011】

さらに、複数の計算機の間でデータ配置に関する管理情報を共有する仕組み有するデータ保管システムでは、単純にデータを移行してしまうと計算機側の管理情報と不整合とな

50

り、データの入出力ができなくなるおそれがある。

【0012】

本発明は、データ保管システムの負荷を増大させることなく、計算機の性能とストレージ装置の負荷との均衡を保ちながら構成を管理する技術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】

【0013】

本発明の代表的な一形態では、クライアントによってデータが格納されるデータ保管サーバと、前記データ保管サーバに記憶領域を提供するストレージ装置と、前記データ保管サーバ及び前記ストレージ装置にアクセス可能な管理計算機と、を含む計算機システムであって、前記データ保管サーバは、前記ストレージ装置に接続される第1インタフェースと、前記第1インタフェースに接続される第1プロセッサと、前記第1プロセッサに接続される第1メモリと、を備え、前記記憶領域が割り当てられたデータ保管領域（データ保管記憶領域）を前記クライアントに提供し、前記ストレージ装置は、前記データ保管サーバに接続される第2インタフェースと、前記第2インタフェースに接続される第2プロセッサと、前記第2プロセッサに接続される第2メモリと、を備え、前記管理計算機は、前記データ保管サーバ及び前記ストレージ装置に接続される第3インタフェースと、前記第3インタフェースに接続される第3プロセッサと、前記第3プロセッサに接続される第3メモリと、を備え、前記データ保管サーバは、前記データと前記データが格納されるデータ保管サーバとの対応を含むデータ配置管理情報、及び前記データ保管領域と前記記憶領域との対応を含むデータ保管記憶領域構成情報を管理し、前記データ配置管理情報は、前記データに基づいて決定されたデータ保管サーバに管理され、前記管理計算機は、前記データ保管サーバの性能を含むデータ保管サーバ性能管理情報、及び前記記憶領域の負荷状態を含む記憶領域稼働情報を管理し、前記データ保管サーバは、前記クライアントからデータの読み出し要求を受け付けた場合には、前記要求されたデータに基づいて、前記要求されたデータに対応するデータ配置管理情報を管理するデータ保管サーバを特定し、前記特定されたデータ保管サーバによって管理されるデータ配置管理情報を取得し、前記取得されたデータ配置管理情報に基づいて、前記要求されたデータを格納するデータ保管サーバを特定し、前記要求されたデータを格納するデータ保管サーバから、前記要求されたデータを取得し、前記管理計算機は、前記データ保管サーバ性能管理情報及び前記記憶領域稼働情報に基づいて、前記記憶領域との対応を変更するデータ保管領域を選択し、前記記憶領域稼働情報に基づいて、前記選択されたデータ保管領域に新たに割り当てられる記憶領域を選択し、前記選択されたデータ保管領域を提供するデータ保管サーバに、前記選択された新たに割り当てられる記憶領域を通知し、前記データ保管サーバは、前記管理計算機から通知された記憶領域を前記データ保管領域に割り当て、前記データ配置管理情報及び前記データ保管記憶領域構成情報を更新する。

10

20

30

【発明の効果】

【0014】

本発明の一形態によれば、計算機の性能とストレージ装置の負荷とのバランスを改善するようにデータ保管システムの構成を管理することができる。

【発明を実施するための最良の形態】

40

【0015】

図1は、本発明の実施の形態のストレージエリアネットワークの構成を表す図である。

【0016】

ストレージエリアネットワークは、データ入出力用ネットワーク400及び管理用ネットワーク600によって構成される。

【0017】

データ入出力用ネットワーク400は、ストレージ装置100、クライアント200、及びデータ保管サーバ300を含む。ストレージ装置100、クライアント200、及びデータ保管サーバ300は、データ入出力用ネットワーク400を介して互いに接続することによって、相互にデータを入出力する。データ入出力用ネットワーク400は、図1

50

では太線で表示されている。データ入出力用ネットワーク４００は、ファイバチャネル又はイーサネット（「イーサネット」は登録商標、以下同じ）などの従来技術によるネットワークである。

【００１８】

管理用ネットワーク６００は、ファイバチャネル又はイーサネットなどの従来技術によるネットワークである。ストレージ装置１００及びデータ保管サーバ３００は、管理用ネットワーク６００を介して管理計算機５００に接続する。

【００１９】

ストレージ装置１００は、磁気記憶装置などの記憶デバイスを搭載し、クライアント２００によって読み書きされるデータの記憶領域を提供する。ストレージ装置１００の構成については、図２にて後述する。クライアント２００は、データ保管サーバ３００を介してストレージ装置１００に格納されたデータを生成及び更新する。

【００２０】

データ保管サーバ３００は、ストレージ装置１００に格納されたデータを管理する。データ保管サーバ３００は、ストレージ装置１００によって提供された記憶領域をマウントし、クライアント２００に記憶領域を提供する。データ保管サーバ３００の構成については、図３にて後述する。データ保管システムは、主として、データ保管サーバ３００及びストレージ装置１００によって構成される。

【００２１】

管理計算機５００は、管理用ネットワーク６００を介してストレージ装置１００及びデータ保管サーバ３００に接続し、データ保管システムを管理する。管理計算機５００の構成については、図４にて後述する。

【００２２】

本発明の実施の形態では、管理用ネットワーク６００と、データ入出力用ネットワーク４００とが、それぞれ独立した形態となっているが、双方の機能を兼ねる単一のネットワークであってもよい。

【００２３】

図２は、本発明の実施の形態のストレージ装置１００の構成を示す図である。

【００２４】

ストレージ装置１００は、データ入出力用通信インタフェース１４０、管理用通信インタフェース１５０、ストレージコントローラ１９０、プログラムメモリ１０００、データ入出力用キャッシュメモリ１６０及び磁気記憶装置１２０を含む。データ入出力用通信インタフェース１４０、管理用通信インタフェース１５０、プログラムメモリ１０００、データ入出力用キャッシュメモリ１６０及び磁気記憶装置１２０は、ストレージコントローラ１９０を介して互いに接続される。

【００２５】

データ入出力用通信インタフェース１４０は、データ入出力用ネットワーク４００を介してデータ保管サーバ３００に接続する。管理用通信インタフェース１５０は、管理用ネットワーク６００を介して管理計算機５００に接続する。なお、データ入出力用通信インタフェース１４０及び管理用通信インタフェース１５０の数は任意である。また、データ入出力用通信インタフェース１４０は、管理用通信インタフェース１５０と独立した構成とする必要はなく、データ入出力用通信インタフェース１４０から管理情報を入出力し、管理用通信インタフェース１５０と共用してもよい。

【００２６】

ストレージコントローラ１９０は、ストレージ装置１００を制御するプロセッサを搭載する。データ入出力用キャッシュメモリ１６０は、記憶領域に対する入出力を高速化するための一時記憶領域である。データ入出力用キャッシュメモリ１６０は、揮発性メモリで構成されることが一般的であるが、不揮発性メモリ又は磁気記憶装置で代用することも可能である。なお、データ入出力用キャッシュメモリ１６０の個数及び容量に制限はない。磁気記憶装置１２０は、クライアント２００によって読み書きされるデータを格納する。

10

20

30

40

50

【0027】

プログラムメモリ1000は、ストレージ装置100で実行される処理に必要なプログラム及び制御情報を格納する。プログラムメモリ1000は、磁気記憶装置又は揮発性半導体メモリによって構成される。プログラムメモリ1000に格納される制御プログラム及び制御情報は、図5にて後述する。

【0028】

図3は、本発明の実施の形態のデータ保管サーバ300の構成を示す図である。

【0029】

データ保管サーバ300は、データ入出力用通信インタフェース340、管理用通信インタフェース350、入力用インタフェース370、出力用インタフェース375、演算処理装置380、磁気記憶装置320及びデータ入出力用キャッシュメモリ360を含む。

10

【0030】

データ入出力用通信インタフェース340、管理用通信インタフェース350、入力用インタフェース370、出力用インタフェース375、演算処理装置380、磁気記憶装置320及びデータ入出力用キャッシュメモリ360は、通信バス390を介して互いに接続される。データ保管サーバ300は、汎用計算機(PC)で実現可能なハードウェア構成である。

【0031】

データ入出力用通信インタフェース340は、データ入出力用ネットワーク400を介してストレージ装置100及びクライアント200に接続し、データを入出力する。管理用通信インタフェース150は、管理用ネットワーク600を介して管理計算機500に接続し、管理情報を入出力する。なお、データ入出力用通信インタフェース340及び管理用通信インタフェース350の数は任意である。また、データ入出力用通信インタフェース340は、管理用通信インタフェース350と独立した構成とする必要はなく、データ入出力用通信インタフェース340から管理情報を入出力し、管理用通信インタフェース350と共用してもよい。

20

【0032】

入力用インタフェース370は、利用者が情報を入力するための機器と接続し、例えば、キーボード及びマウスなどに接続する。出力用インタフェース375は、利用者に情報を出力するための機器と接続し、例えば、汎用ディスプレイなどに接続する。演算処理装置380は、各種演算を実行し、CPU又はプロセッサに相当する。磁気記憶装置320は、オペレーティングシステム及びアプリケーションなどのソフトウェアを格納する。

30

【0033】

データ入出力用キャッシュメモリ360は、揮発性メモリなどによって構成され、磁気記憶装置320又はストレージ装置100によって提供される記憶領域に対するデータ入出力を高速化するために使用される。データ入出力用キャッシュメモリ360は、揮発性メモリによる実装が一般的であるが、不揮発性メモリ又は磁気記憶装置で構成してもよい。なお、データ入出力用キャッシュメモリ360の個数及び容量に制限はない。

【0034】

プログラムメモリ3000は、データ保管サーバ300で実行される処理に必要なプログラム及び制御情報を格納する。プログラムメモリ3000は、磁気記憶装置又は揮発性半導体メモリによって構成される。プログラムメモリ3000に格納されるプログラム及び制御情報は、図6にて後述する。

40

【0035】

図4は、本発明の実施の形態の管理計算機500の構成を示す図である。

【0036】

管理計算機500は、管理用通信インタフェース550、入力用インタフェース570、出力用インタフェース575、演算処理装置580、磁気記憶装置520、プログラムメモリ5000及びデータ入出力用キャッシュメモリ560を含む。

50

【 0 0 3 7 】

管理用通信インタフェース 5 5 0、入力用インタフェース 5 7 0、出力用インタフェース 5 7 5、演算処理装置 5 8 0、磁気記憶装置 5 2 0、プログラムメモリ 5 0 0 0 及びデータ入出力用キャッシュメモリ 5 6 0 は、通信バス 5 9 0 を介して互いに接続される。管理計算機 5 0 0 は、汎用計算機 (P C) で実現可能なハードウェア構成となっており、各部の機能は、図 3 に示されたデータ保管サーバ 3 0 0 と同じである。

【 0 0 3 8 】

プログラムメモリ 5 0 0 0 は、管理計算機 5 0 0 で実行される処理に必要なプログラム及び情報を格納する。プログラムメモリ 5 0 0 0 は、磁気記憶装置又は揮発性半導体メモリによって構成される。プログラムメモリ 5 0 0 0 に格納されるプログラム及び情報は、
10 図 7 にて後述する。

【 0 0 3 9 】

図 5 は、本発明の実施の形態のストレージ装置 1 0 0 のプログラムメモリ 1 0 0 0 に格納される制御プログラム及び制御情報の一例を示す図である。

【 0 0 4 0 】

プログラムメモリ 1 0 0 0 には、記憶領域構成管理プログラム 1 0 0 1、RAID グループ構成情報 1 0 0 2、記憶領域構成情報 1 0 0 3、論理記憶ユニット構成情報 1 0 0 4、記憶領域稼働監視プログラム 1 0 0 5 及び記憶領域稼働情報 1 0 0 6 が含まれる。

【 0 0 4 1 】

記憶領域構成管理プログラム 1 0 0 1 は、ストレージコントローラ 1 9 0 に搭載されたプロセッサに実行されることによって、後述する記憶領域構成情報 1 0 0 3 に基づいて、
20 データ保管サーバ 3 0 0 に提供する記憶領域を管理及び制御する。

【 0 0 4 2 】

RAID グループ構成情報 1 0 0 2 は、磁気記憶装置 1 2 0 の集合による RAID グループの構成情報である。RAID グループ構成情報 1 0 0 2 の詳細については、図 9 にて後述する。

【 0 0 4 3 】

記憶領域構成情報 1 0 0 3 は、RAID グループが論理的な単位に分割された記憶資源の単位である記憶領域の構成情報である。記憶領域構成情報 1 0 0 3 の詳細については、
30 図 1 0 にて後述する。

【 0 0 4 4 】

論理記憶ユニット構成情報 1 0 0 4 は、データ保管サーバ 3 0 0 に提供する記憶資源の単位である論理記憶ユニットの構成情報である。論理記憶ユニット構成情報 1 0 0 4 の詳細については、図 1 1 にて後述する。

【 0 0 4 5 】

記憶領域稼働監視プログラム 1 0 0 5 は、ストレージコントローラ 1 9 0 に搭載されたプロセッサに実行されることによって、記憶領域ごとの稼働情報を監視し、記憶領域稼働情報 1 0 0 6 に稼働情報を記録する。記憶領域稼働情報 1 0 0 6 は、記憶領域ごとの入出力命令数、読み出し又は書き込みデータ量、ハードディスク稼働率、応答時間などの稼働情報が記録される。記憶領域稼働情報 1 0 0 6 の詳細については、図 1 2 にて後述する。
40

【 0 0 4 6 】

図 6 は、本発明の実施の形態のデータ保管サーバ 3 0 0 のプログラムメモリ 3 0 0 0 に格納される制御プログラム及び制御情報の一例を示す図である。

【 0 0 4 7 】

プログラムメモリ 3 0 0 0 には、データ入出力プログラム 3 0 0 1、記憶領域構成管理プログラム 3 0 0 2、データ保管記憶領域構成情報 3 0 0 3、データ配置管理プログラム 3 0 0 4、データ配置管理情報格納サーバ計算プログラム 3 0 0 5、データ配置管理情報 3 0 0 6 及びハッシュ計算情報 3 0 0 7 が含まれる。

【 0 0 4 8 】

データ入出力プログラム 3 0 0 1 は、データ保管サーバ 3 0 0 の演算処理装置 3 8 0 に
50

実行されることによって、クライアント 200 からの要求に基づいて、ストレージ装置 100 に格納されたデータを読み書きし、クライアント 200 に処理結果を応答する。

【0049】

データ保管サーバ 300 は、前述したように、ストレージ装置 100 によって提供された記憶領域（論理記憶ユニット）がマウントされ、マウントされた記憶領域をクライアント 200 に提供する。記憶領域構成管理プログラム 3002 は、マウントされた記憶領域を管理する。

【0050】

データ保管記憶領域構成情報 3003 は、データ保管サーバ 300 にマウントされたデータ保管記憶領域（データ保管領域）の構成情報を格納する。具体的には、データ保管記憶領域と論理記憶ユニットとの対応が格納される。データ保管記憶領域構成情報 3003 の詳細については、図 13 にて後述する。

10

【0051】

データ配置管理プログラム 3004 は、データ保管サーバ 300 の演算処理装置 380 に実行されることによって、データ配置管理情報格納サーバ計算プログラム 3005 及びデータ配置管理情報 3006 に基づいて、クライアント 200 から要求されたデータの格納場所を特定する。また、新たにデータを保管する場合には、データを格納する記憶領域を選定する。

【0052】

データ配置管理情報格納サーバ計算プログラム 3005 は、データ保管サーバ 300 の演算処理装置 380 に実行されることによって、要求されたデータの格納場所を記録したデータ配置管理情報 3006 を格納するデータ保管サーバ 300 を特定する。例えば、要求されたデータに基づいてハッシュ計算を行い、算出されたハッシュ値及びハッシュ計算情報 3007 に基づいてデータ保管サーバ 300 を特定する。具体的な計算手順については、図 19 に示すファイルの保管場所を特定する処理にて説明する。

20

【0053】

データ配置管理情報 3006 は、ファイル名、当該ファイルが格納されているデータ保管サーバ 300 及びデータ保管記憶領域の対応を格納する。データ配置管理情報 3006 の詳細については、図 14 にて後述する。

【0054】

ハッシュ計算情報 3007 は、データ配置管理情報格納サーバ計算プログラム 3005 によって、データ配置管理情報 3006 を格納するデータ保管サーバ 300 を特定するために使用され、ハッシュ値とデータ保管サーバ 300 との対応を格納する。ハッシュ計算情報 3007 の詳細については、図 15 にて後述する。

30

【0055】

図 7 は、本発明の実施の形態の管理計算機 500 のプログラムメモリ 5000 に格納される制御プログラム及び制御情報の一例を示す図である。

【0056】

ストレージ装置論理記憶ユニット構成情報 5001 は、ストレージ装置 100 に格納された論理記憶ユニット構成情報 1004 に対応する。データ保管サーバデータ保管記憶領域構成情報 5002 は、データ保管サーバ 300 に格納されたデータ保管記憶領域構成情報 3003 に対応する。

40

【0057】

構成情報更新サービスプログラム 5004 は、管理計算機 500 の演算処理装置 580 に実行されることによって、ストレージ装置 100 及びデータ保管サーバ 300 と定期的に通信し、ストレージ装置論理記憶ユニット構成情報 5001 及びデータ保管サーバデータ保管記憶領域構成情報 5002 を最新の状態に維持する。

【0058】

ストレージ装置論理記憶ユニット稼働情報 5005 は、ストレージ装置 100 に格納された記憶領域稼働情報 1006 を、論理記憶ユニットごとに集計した稼働情報である。ス

50

ストレージ装置論理記憶ユニット稼働情報 5 0 0 5 の詳細については、図 1 6 にて後述する。

【 0 0 5 9 】

稼働情報更新サービスプログラム 5 0 0 6 は、管理計算機 5 0 0 の演算処理装置 5 8 0 に実行されることによって、ストレージ装置 1 0 0 と定期的に通信し、ストレージ装置論理記憶ユニット稼働情報 5 0 0 5 を最新の状態に維持する。

【 0 0 6 0 】

データ保管サーバ性能情報 5 0 0 7 は、データ保管サーバ 3 0 0 の性能に基づいて、管理者によって定義された順位を含む情報である。データ保管サーバ 3 0 0 の性能とは、例えば、演算処理装置 3 8 0 の処理性能であってもよいし、データ入出力用通信インタフェース 3 4 0 の処理性能であってもよい。ストレージ装置論理記憶ユニット稼働情報 5 0 0 5 の詳細については、図 1 7 にて後述する。

【 0 0 6 1 】

データ保管サーバ性能情報更新サービスプログラム 5 0 0 8 は、管理計算機 5 0 0 の演算処理装置 5 8 0 に実行されることによって、データ保管サーバ性能情報 5 0 0 7 を更新又は管理する。

【 0 0 6 2 】

データ保管記憶領域配置計算プログラム 5 0 0 9 は、管理計算機 5 0 0 の演算処理装置 5 8 0 に実行されることによって、データ保管サーバ 3 0 0 の性能と論理記憶ユニットの負荷のバランスが最適になるようにデータ保管記憶領域の配置を計算する。例えば、最も性能の高いデータ保管サーバ 3 0 0 に最も負荷の高いデータ保管記憶領域を対応付ける。

【 0 0 6 3 】

データ配置管理情報更新プログラム 5 0 1 0 は、管理計算機 5 0 0 の演算処理装置 5 8 0 に実行されることによって、データ保管サーバ 3 0 0 にデータ配置管理情報 3 0 0 6 の更新を指示する。

【 0 0 6 4 】

以下、データ保管システムの構成を、図 8 を参照しながら説明し、図 8 に示したデータ保管システムの構成情報を図 9 から図 1 7 を参照しながら説明する。

【 0 0 6 5 】

図 8 は、本発明の実施の形態のデータ保管システムの構成の一例を示す図である。

【 0 0 6 6 】

図 8 に示すデータ保管システムには、データ保管サーバ 3 0 0 A、3 0 0 B 及び 3 0 0 C が含まれ、さらに、ストレージ装置 1 0 0 A、1 0 0 B 及び 1 0 0 C が含まれる。なお、図 8 に示すようにデータ保管システムに複数のデータ保管サーバ 3 0 0 及びストレージ装置 1 0 0 がそれぞれ含まれ、これらの機器を個別に識別する必要がある場合には、数字による符号に加えアルファベットを連結した符号を付与する。例えば、データ保管サーバ 3 0 0 については、符号を 3 0 0 A、3 0 0 B 及び 3 0 0 C とする。また、共通する構成又は処理については、アルファベットが付与されていない符号（例えば、データ保管サーバ 3 0 0）で説明する。

【 0 0 6 7 】

データ保管サーバ 3 0 0 A のデータ入出力用通信インタフェース 3 4 0 A は、データ入出力用ネットワーク 4 0 0 を介してストレージ装置 1 0 0 A のデータ入出力用通信インタフェース 1 4 0 A に接続される。同様に、データ保管サーバ 3 0 0 B はストレージ装置 1 0 0 B に接続され、データ保管サーバ 3 0 0 C はストレージ装置 1 0 0 C に接続される。

【 0 0 6 8 】

ストレージ装置 1 0 0 A のデータ入出力用通信インタフェース 1 4 0 A は、識別情報 “ 5 0 : 0 0 : 0 1 : 1 E : 0 A : E 8 : 0 2 ” で識別される。ストレージ装置 1 0 0 A のデータ入出力用通信インタフェース 1 4 0 A には、“ L U - 1 1 ” で識別される論理記憶ユニット 1 1 0 0 A 及び “ L U - 1 2 ” で識別される論理記憶ユニット 1 1 0 0 B が登録されている。さらに、論理記憶ユニット 1 1 0 0 A は “ L D - 0 1 ” で識別される論理記

10

20

30

40

50

憶領域で構成され、論理記憶ユニット 1100B は“LD-02”で識別される論理記憶領域で構成される。以上の構成は、図 11 に示す論理記憶ユニット構成情報 1004 において定義されている。

【0069】

さらに、論理記憶ユニット構成情報 1004 には、ストレージ装置 100 のデータ入出力用通信インタフェース 140 と、アクセス可能なデータ保管サーバ 300 のデータ入出力用通信インタフェース 340 との対応が格納されている。具体的には、ストレージ装置 100A のデータ入出力用通信インタフェース 140A には、データ保管サーバ 300A のデータ入出力用通信インタフェース 340A (識別情報“50:00:01:1E:0A:E8:A1”)が対応付けられている。したがって、論理記憶ユニット 1100A 及び 1100B は、“192.168.0.1”で識別されるデータ保管サーバ 300A によってアクセス可能となる。

10

【0070】

データ保管サーバ 300A は、論理記憶ユニット 1100A を、“/mount/data1”で識別されるデータ保管記憶領域 3100A にマウントしてクライアント 200 に提供する。論理記憶ユニット 1100B についても同様に、“/mount/data2”で識別されるデータ保管記憶領域 3100B にマウントされる。論理記憶ユニットとデータ保管記憶領域との関係は、図 13 に示すデータ保管記憶領域構成情報 3003 に定義される。

20

【0071】

データ保管サーバ 300B 及び 300C のデータ保管記憶領域の構成についても同様に定義される。図 8 に示したデータ保管システムを構成するための構成情報を、以下、説明する。

【0072】

図 9 は、本発明の実施の形態のストレージ装置 100 に格納される RAID グループ構成情報 1002 の一例を示す図である。

【0073】

RAID グループ構成情報 1002 は、RAID グループと RAID グループを構成する磁気記憶装置との対応関係を格納する。RAID グループ構成情報 1002 は、RAID グループ識別情報 10021 及び磁気記憶装置識別情報 10022 を含む。

30

【0074】

RAID グループ識別情報 10021 は、ストレージ装置 100 に備えられる RAID グループを一意に識別する識別子である。

【0075】

磁気記憶装置識別情報 10022 は、RAID グループ識別情報 10021 によって特定される RAID グループを構成する磁気記憶装置 120 を一意に識別する識別子である。例えば、RAID グループ「RG-01」は、磁気記憶装置「HD-01」、「HD-02」、「HD-03」及び「HD-04」によって構成される。

【0076】

図 10 は、本発明の実施の形態のストレージ装置 100 に格納される記憶領域構成情報 1003 の一例を示す図である。

40

【0077】

記憶領域構成情報 1003 は、記憶領域識別情報 10031、RAID グループ識別情報 10032、開始ブロックアドレス 10033 及び終了ブロックアドレス 10034 を含む。

【0078】

記憶領域識別情報 10031 は、記憶領域を識別する識別子である。RAID グループ識別情報 10032 は、RAID グループを識別する識別子である。記憶領域識別情報 10031 によって識別される記憶領域は、RAID グループ識別情報 10032 によって識別される RAID グループに定義された論理的な記憶領域である。

50

【0079】

開始ブロックアドレス10033は、記憶領域識別情報10031によって識別される記憶領域が格納される物理領域の開始ブロックアドレスである。一方、終了ブロックアドレス10034は、記憶領域識別情報10031によって識別される記憶領域が格納される物理領域の終了ブロックアドレスである。

【0080】

図11は、本発明の実施の形態のストレージ装置100に格納される論理記憶ユニット構成情報1004の一例を示す図である。

【0081】

論理記憶ユニット構成情報1004は、ストレージ装置100のデータ入出力用通信インタフェース140と、記憶領域と、さらに、外部からアクセスが許可されているインタフェース識別情報との対応が定義される。

10

【0082】

論理記憶ユニット構成情報1004は、通信インタフェース識別情報10041、記憶ユニット識別情報10042、記憶領域識別情報10043及びアクセス許可通信インタフェース識別情報10044を含む。

【0083】

通信インタフェース識別情報10041は、データ入出力用通信インタフェース140を一意に識別する識別子である。通信インタフェース識別情報10041には、例えば、WWN (World Wide Name) が格納される。

20

【0084】

記憶ユニット識別情報10042は、論理記憶ユニットを一意に識別する識別子である。論理記憶ユニットは、ストレージ装置100に接続されたデータ保管サーバ300によってアクセス可能な記憶資源の単位である。データ保管サーバ300は、論理記憶ユニットをマウントし、データ保管記憶領域としてクライアント200に提供する。

【0085】

記憶領域識別情報10043は、ストレージ装置100によって提供される論理的な記憶領域を一意に識別する識別子である。

【0086】

アクセス許可通信インタフェース識別情報10044は、記憶ユニット識別情報10042によって識別される記憶ユニットにアクセスが許可された機器の通信インタフェースを識別する情報である。アクセス許可通信インタフェース識別情報10044を定義することによって、特定のデータ保管サーバ300からのみアクセスを可能とし、セキュリティを向上させることができる。

30

【0087】

図12は、本発明の実施の形態のストレージ装置100に格納される記憶領域稼働情報1006の一例を示す図である。

【0088】

記憶領域稼働情報1006は、時刻10061ごとの稼働状態10062が格納される。また、記憶領域稼働情報1006は、記憶領域ごとの稼働監視結果が時系列順に記録される。

40

【0089】

時刻10061は、稼働状況が記録された時刻情報である。稼働状態10062は、具体的には、入出力要求数 (I/O数)、読み書きされたデータ量 (GB又はMB/sec)、ハードディスク稼働率 (%)、応答時間 (ms) などの性能指標である。図12に示す記憶領域稼働情報1006では、所定期間における入出力要求数 (I/O数) である。

【0090】

図13は、本発明の実施の形態のデータ保管サーバ300に格納されるデータ保管記憶領域構成情報3003の一例を示す図である。

【0091】

50

データ保管記憶領域構成情報 3003 は、データ保管記憶領域識別情報 30031、通信インタフェース識別情報 30032、記憶ユニット識別情報 30033、容量使用率 30034 及び転送データ量 30035 を含む。

【0092】

データ保管記憶領域識別情報 30031 は、データ保管サーバ 300 によって運用されるファイルシステムのマウントポイントに該当する。データ保管記憶領域識別情報 30031 は、ストレージ装置 100 によって提供される論理記憶ユニットをファイルシステム上で参照するための識別情報である。

【0093】

通信インタフェース識別情報 30032 は、ストレージ装置 100 のデータ入出力用通信インタフェース 140 を一意に識別する識別子である。 10

【0094】

記憶ユニット識別情報 30033 は、通信インタフェース識別情報 30032 によって識別されるデータ入出力用通信インタフェース 140 に登録された論理記憶ユニットを一意に識別するための識別子である。すなわち、データ保管サーバ 300 によって提供されるデータ保管記憶領域の実体は、データ入出力用ネットワーク 400 を介して接続され、ストレージ装置 100 のデータ入出力用通信インタフェース 140 に登録された論理記憶ユニットである。

【0095】

容量使用率 30034 は、データ保管記憶領域の容量に対して、記録されているデータ量が占める容量の割合である。転送データ量 30035 は、データ保管記憶領域において読み書きされたデータ量である。容量使用率 30034 及び転送データ量 30035 は、新たにデータが格納されるデータ保管記憶領域を選択する場合などに参照される。 20

【0096】

図 14 は、本発明の実施の形態のデータ保管サーバ 300 に格納されるデータ配置管理情報 3006 の一例を示す図である。

【0097】

データ配置管理情報 3006 は、ファイル識別情報 30061、データ保管サーバ識別情報 30062 及びデータ保管記憶領域識別情報 30063 を含む。

【0098】

ファイル識別情報 30061 は、保管されるデータに対応するファイルを識別する識別情報である。具体的にはファイル名である。 30

【0099】

データ保管サーバ識別情報 30062 は、当該データが保管されているデータ保管サーバ 300 の識別情報である。本発明の実施の形態では、データ保管サーバ 300 の通信インタフェースの IP アドレスによってデータ保管サーバ 300 を識別する。

【0100】

データ保管記憶領域識別情報 30063 は、当該データが格納されているデータ保管記憶領域を一意に識別するための識別子である。

【0101】

したがって、データ保管システムによって管理されるデータは、データ保管サーバ識別情報 30062 によって識別されるデータ保管サーバ 300 の、データ保管記憶領域識別情報 30063 によって識別されるデータ保管記憶領域に格納された、ファイル識別情報 30061 によって識別されるファイルに該当する。 40

【0102】

図 15 は、本発明の実施の形態のデータ保管サーバ 300 に格納されるハッシュ計算情報 3007 の一例を示す図である。

【0103】

ハッシュ計算情報 3007 は、サーバ識別情報 30071 及びサーバ検索情報 30072 を含む。 50

【0104】

サーバ識別情報30071は、読み出し又は保管が要求されたデータ（ファイル）の格納場所を記録したデータ配置管理情報3006を格納するデータ保管サーバ300（データ配置管理情報格納サーバ）の識別情報である。サーバ検索情報30072は、データ保管サーバ300に対応するハッシュ値である。

【0105】

データ保管サーバ300は、ファイルの読み出し又は保管が要求されると、データ配置管理情報格納サーバ計算プログラム3005によって、要求されたファイルのファイル名に基づいてハッシュ値を計算する。続いて、ハッシュ計算情報3007を参照して、算出されたハッシュ値と一致するサーバ検索情報30072に対応するサーバ識別情報30071を取得し、データ配置管理情報格納サーバを特定する。

10

【0106】

図16は、本発明の実施の形態の管理計算機500に格納されるストレージ装置論理記憶ユニット稼働情報5005の一例を示す図である。

【0107】

ストレージ装置論理記憶ユニット稼働情報5005は、ストレージ装置100で収集された記憶領域稼働情報1006を集約することによって取得される。具体的には、管理計算機500で稼働情報更新サービスプログラム5006が実行されることによって、各ストレージ装置100から記憶領域稼働情報1006を収集する。

20

【0108】

ストレージ装置論理記憶ユニット稼働情報5005は、時刻50051ごとに記憶領域ごとに稼働状態50052が格納される。

【0109】

時刻50051は、稼働状況が記録された時刻情報である。稼働状態50052は、具体的には、入出力要求数（I/O数）、読み書きされたデータ量（GB又はMB/sec）、ハードディスク稼働率（%）、応答時間（ms）などの性能指標である。図16に示すストレージ装置論理記憶ユニット稼働情報5005では、所定期間の入出力要求数（I/O数）である。

【0110】

図17は、本発明の実施の形態の管理計算機500に格納されるデータ保管サーバ性能情報5007の一例を示す図である。

30

【0111】

データ保管サーバ性能情報5007は、データ保管サーバ300の性能情報である。データ保管サーバ性能情報5007は、データ入出力性能ランク50071、データ保管サーバ識別情報50072及び通信インタフェース識別情報50073を含む。

【0112】

データ入出力性能ランク50071は、データ保管サーバ300の性能を示す指標である。データ保管サーバ識別情報50072は、データ保管サーバ300の識別情報である。通信インタフェース識別情報50073は、データ保管サーバ300に搭載された通信インタフェースの識別情報である。

40

【0113】

図17に示すデータ保管サーバ性能情報5007では、データ保管サーバ300に搭載されたデータ入出力用通信インタフェース340ごとに、データ入出力性能ランク50071に示された順位によって性能をランク付けしている。

【0114】

また、本発明の実施の形態では、データ保管サーバ性能情報5007が管理者によって定義されることを想定しているが、データ保管サーバ300の演算処理装置380又はデータ入出力用通信インタフェース340の性能などに基づいて、動的にランク付けするようにしてもよい。

【0115】

50

以上、本発明の実施の形態の構成について説明した。以下、本発明の実施の形態における処理について説明する。まず、図18を参照しながら概要を説明し、図19から図31までのフローチャートを参照しながら詳細を説明する。

【0116】

図18は、本発明の実施の形態のデータ保管システムの性能を管理するための手順の概要を示す図である。

【0117】

図18に示したデータ保管システムは、管理計算機500、三台のデータ保管サーバ(300A~300C)及びストレージ装置100を含む。データ保管サーバ300A及び300Bは相対的に性能が低く、データ保管サーバ300Cは相対的に性能が高い。また、ストレージ装置100は、ボリューム1~4を提供する。各ボリュームは、論理記憶ユニットに対応する。

10

【0118】

まず、指定されたファイルをデータ保管システムから読み出す処理について説明する。

【0119】

本発明の実施の形態のデータ保管システムからファイルを読み出す場合、クライアント200は、任意のデータ保管サーバ300にファイルの読み出し要求を送信すればよい。一例として、図18に示すように、クライアント200からデータ保管サーバ300Cにファイル“F12”の読み出し要求を送信した場合について説明する。

【0120】

データ保管サーバ300Cは、まず、読み出しを要求されたファイル“F12”の格納場所が管理されているデータ配置管理情報格納サーバを特定する。データ配置管理情報格納サーバは、例えば、ファイル名に基づいて一意に決定されるようにあらかじめ定められている。具体的な特定方法については後述する。ここで、データ配置管理情報格納サーバをデータ保管サーバ300Aとする。

20

【0121】

データ保管サーバ300Aのデータ配置管理情報3006Aを参照すると、ファイル“F12”がデータ保管サーバ300Cによって管理されていることがわかる。したがって、データ保管サーバ300Cは、自身にマウントされているボリューム3からファイル“F12”を読み出し、クライアント200に送信する。

30

【0122】

以上のように、任意のデータ保管サーバ300にファイルの読み出しを要求することによって、ファイルの保管場所を意識せずに要求したファイルを読み出すことができる。

【0123】

続いて、データ保管サーバ300の性能及びストレージ装置100によって提供されるボリュームの負荷状態に基づいて、データ保管システムの構成を変更する手順について説明する。図18には構成変更後の状態が示されており、構成変更前には、データ保管サーバ300Bとボリューム3とが接続され、データ保管サーバ300Cとボリューム4とが接続されていたものとする。以上の状態で本発明の構成変更手順を適用する手順について説明する。

40

【0124】

まず、管理計算機500は、各ボリュームの稼働状況(負荷状態)を監視する(ステップS1)。負荷状態は、前述したように、稼働情報更新サービスプログラム5006を実行することによって、ストレージ装置論理記憶ユニット稼働情報5005に記録される。

【0125】

図18に示す状態では、ボリューム3の負荷が大きく、ボリューム4の負荷が小さいことがわかる。なお、各ボリュームの稼働状況は各ボリュームの下部に示しており、単位時間あたりの転送データ量としている。

【0126】

しかし、構成変更前の状態では、低性能のデータ保管サーバ300Bに高負荷のボリュ

50

ーム3に接続されているため、データ保管サーバ300Bの処理が追いつかず、システム全体の性能が低下していることが考えられる。一方、高性能のデータ保管サーバ300Cに低負荷のボリューム4が接続されているため、データ保管サーバ300Cの性能が十分に発揮されていないと考えられる。

【0127】

そこで、低性能のデータ保管サーバ300Bに低負荷のボリューム4を接続し、高性能のデータ保管サーバ300Cに高負荷のボリューム3を接続することによって、性能と負荷のバランスを改善させる。構成を変更するデータ保管サーバ及びボリュームの組合せは、データ保管サーバの性能とボリュームの稼働状況に基づいて決定される。例えば、処理性能に対して負荷の高いボリュームに接続されているデータ保管サーバと、処理性能に対して負荷の低いボリュームが接続されているデータ保管サーバとを選択し、接続を入れ替える。

10

【0128】

構成を変更する手順を説明すると、まず、管理計算機500は、データ配置管理プログラム5003を実行することによって、ボリュームに新たに接続されるデータ保管サーバ300によるアクセスを許可するようにストレージ装置に要求する(ステップS2)。具体的には、データ保管サーバ300Cによるボリューム3へのアクセス及びデータ保管サーバ300Bによるボリューム4へのアクセスを許可するように、ストレージ装置100に要求する。

【0129】

続いて、管理計算機500は、データ配置管理プログラム5003を実行することによって、データ保管サーバ300に変更先のボリュームをマウントするように指示する(ステップS3)。具体的には、データ保管サーバ300Bがボリューム3をマウントし、データ保管サーバ300Cがボリューム4をマウントするように指示する。なお、ボリュームをマウントする機能は、データ保管サーバ300で実行されるオペレーションシステム(OS)によって提供される。

20

【0130】

最後に、管理計算機500は、データ保管サーバ300の構成情報を変更後の状態に更新する(ステップS4)。具体的には、データ配置管理情報更新プログラム5010を実行することによって、データ配置管理情報3006の更新を指示する。図18を参照すると、データ保管サーバ300Aのデータ配置管理情報3006Aでは、ファイル“F12”を管理するサーバがデータ保管サーバ300Bからデータ保管サーバ300Cに変更されている。同様に、データ保管サーバ300Cのデータ配置管理情報3006Cでは、ファイル“F31”を管理するサーバがデータ保管サーバ300Cからデータ保管サーバ300Bに変更されている。

30

【0131】

以上の処理によって、データ保管システムの構成を変更することができる。さらに、構成情報を更新する詳細な処理手順については、図25から図31を参照しながら説明する。ここで、構成情報の更新処理について説明する前に、データの保管処理(図19及び図20)、データの読み出し処理(図21及び図22)、構成情報の更新処理(図23)、及び稼働情報の更新処理(図24)について図面を参照しながら説明する。

40

【0132】

図19及び図20は、本発明の実施の形態のデータ保管システムによるデータ保管処理の手順を示すフローチャートである。図19に示す手順は、主に、データの保管場所を特定し、データ配置管理情報を更新する手順である。図20に示す手順は、主に、実際にファイルを保管する手順である。

【0133】

クライアント200は、データ保管システムにファイルデータを保管する場合には、データ保管システムに含まれる任意のデータ保管サーバ300にファイル保管要求メッセージを送信する(ステップS101)。ファイル保管要求メッセージには、保管されるファ

50

イルデータが含まれる。以後、要求送信先となったデータ保管サーバ300を、データ保管サーバAとする。

【0134】

データ保管サーバAの演算処理装置380は、ファイル保管要求メッセージを受信すると(ステップS102)、データ配置管理プログラム3004を実行する。そして、データ配置管理情報格納サーバ計算プログラム3005によって、当該ファイルデータの格納場所が定義されたデータ配置管理情報3006を格納するデータ配置管理情報格納サーバを特定する(ステップS103)。

【0135】

ここで、データ配置管理情報格納サーバの計算方法について説明する。単純な計算方法としては、保管されるファイルデータに基づいてハッシュ計算を実行する。例えば、ファイル名又はデータの一部を入力値としてハッシュ計算を実行することによって、一意に決定される数値を算出する。算出された数値をデータ保管サーバ300の台数で除算した余りは、データ保管サーバ300の台数未満の数値となる。各データ保管サーバ300に0からデータ保管サーバ300の台数-1の値をあらかじめ対応づけておけば、データ配置管理情報格納サーバ計算プログラム3005による計算結果によって、いずれかのデータ保管サーバ300に対応させることができる。

10

【0136】

さらに具体的に説明すると、保管されるデータのファイル名が“0016.dat”であって、このファイル名に基づいて算出されたハッシュ計算値が“3728”であったとする。算出されたハッシュ値をデータ保管サーバ300の台数(3台)で除算した余りは“2”となる。そこで、図15に示したハッシュ計算情報3007を参照すると、ハッシュ計算値“2”に対応するデータ配置管理情報格納サーバは、識別情報“192.168.0.3”に対応するデータ保管サーバ300となる。

20

【0137】

なお、以上のデータ配置管理情報格納サーバを算出するための計算結果は、いずれのデータ保管サーバ300においても同一の結果が得られる。したがって、クライアント200がいずれのデータ保管サーバ300に入出力要求を行っても、データ配置管理情報格納サーバは一意に特定されるため、クライアント200は任意のデータ保管サーバ300に処理を要求すればよい。

30

【0138】

次に、データ保管サーバAの演算処理装置380は、データ配置管理プログラム3004を実行することによって、データを格納及び保管するデータ保管サーバ300及びデータ保管記憶領域3100を選定する(ステップS104)。選定にあたっては、例えば、データ保管記憶領域構成情報3003を参照し、容量使用率30034に基づいて容量使用率の低いデータ保管記憶領域3100を選択するといった基準を採用すればよい。また、容量ではなく、データ保管記憶領域3100の負荷に応じて選択するようにしてもよい。例えば、データ保管記憶領域構成情報3003の転送データ量30035に基づいて、転送データ量の少ないデータ保管記憶領域3100を選択してもよい。さらに、ストレージ装置論理記憶ユニット稼働情報5005に基づいて、単位時間あたりのディスクアクセス頻度の少ないデータ保管記憶領域3100を選択してもよい。

40

【0139】

次に、データ保管サーバAの演算処理装置380は、データ配置管理プログラム3004によって、ステップS103の処理で特定されたデータ配置管理情報格納サーバに対し、ステップS102で要求された“ファイル名”と、ステップS104で選定されたデータ保管サーバ300及びデータ保管記憶領域3100を記録したデータ配置管理情報更新要求を送信する(ステップS105)。以降、送信先となったデータ配置管理情報格納サーバをデータ保管サーバBとする。

【0140】

データ保管サーバBの演算処理装置380は、データ配置管理情報更新要求を受信する

50

と(ステップS106)、データ配置管理プログラム3004によって、データ配置管理情報3006に新たなエントリを追加する(ステップS107)。ステップS107の処理が完了すると、データ保管サーバAに完了通知を送信する(ステップS108)。

【0141】

データ保管サーバAの演算処理装置380は、完了通知を受信すると(ステップS109)、データ配置管理プログラム3004によって、ステップS104の処理で選定されたデータ保管記憶領域3100を管理するデータ保管サーバ300にファイル格納要求メッセージを送信する(ステップ110)。以降、要求先データ保管サーバをデータ保管サーバCとする。

【0142】

データ保管サーバCの演算処理装置380は、ファイル格納要求メッセージを受信すると(ステップS111)、データ入出力プログラム3001によって、要求されたデータ保管記憶領域3100に受信したデータを格納する(ステップS112)。データの格納が完了すると、完了通知をデータ保管サーバAに送信する(ステップS113)。

【0143】

データ保管サーバAの演算処理装置380は、データ保管サーバCから送信された完了通知を受信すると(ステップS114)、要求元のクライアントに対して完了通知を送信する(ステップS115)。クライアント200が完了通知を受信すると(ステップS116)、データ保管処理は完了する。

【0144】

図21及び図22は、本発明の実施の形態のデータ保管システムに保管されたデータの読み出し処理の手順を示すフローチャートである。図21に示す手順は、主に、データの保管場所を特定する手順である。図22に示す手順は、主に、実際にファイルを読み出す手順である。

【0145】

クライアント200は、まず、任意のデータ保管サーバ300に対してファイル名を指定して読み出し要求メッセージを送信する(ステップS201)。

【0146】

データ保管サーバ300の演算処理装置380は、クライアント200から送信された読み出し要求メッセージを受信し(ステップS202)、データ配置管理プログラム3004を実行する。以降、読み出し要求メッセージを受信したデータ保管サーバ300をデータ保管サーバAとする。

【0147】

データ保管サーバAの演算処理装置380は、データ配置管理情報格納サーバ計算プログラム3005によって、読み出しを要求されたファイルのデータ配置管理情報3006を格納するデータ配置管理情報格納サーバを特定する。算出方法については、図19のステップS103と同様である。

【0148】

データ保管サーバAの演算処理装置380は、データ配置管理プログラム3004によって、特定されたデータ配置管理情報格納サーバにファイル名を含むデータ配置管理情報参照要求メッセージを送信する(ステップS204)。以降、特定されたデータ配置管理情報格納サーバをデータ保管サーバBとする。

【0149】

データ保管サーバBの演算処理装置380は、データ保管サーバAからデータ配置管理情報参照要求メッセージを受信すると(ステップS205)、受信したメッセージに含まれるファイル名にファイル識別情報30061が一致するエントリをデータ配置管理情報3006から検索する(ステップS206)。そして、検索結果をデータ保管サーバAに送信する(ステップS207)。

【0150】

データ保管サーバAの演算処理装置380は、データ保管サーバBからデータ配置管理

10

20

30

40

50

情報 3006 の検索結果を受信すると (ステップ S208)、ステップ S208 の処理で受信したデータ配置管理情報 3006 の検索結果に基づいて、データが保管されているデータ保管サーバ 300 にファイル読み出し要求を送信する (ステップ S209)。具体的には、ファイル読み出し要求が送信されたデータ保管サーバ 300 は、データ保管サーバ識別情報 30062 に記録されている。さらに、ファイル読み出し要求には、データ保管記憶領域識別情報 30063 及びファイル識別情報 30061 が含まれている。以降、データが保管されているデータ保管サーバ 300 をデータ保管サーバ C とする。

【0151】

データ保管サーバ C の演算処理装置 380 は、データ保管サーバ A からファイル読み出し要求を受信すると (ステップ S210)、データ入出力プログラム 3001 を実行することによって、データ保管記憶領域 3100 から要求されたファイルを読み出す (ステップ S211)。さらに、読み出されたファイルをデータ保管サーバ A に送信する (ステップ S212)。

10

【0152】

データ保管サーバ A の演算処理装置 380 は、データ保管サーバ C から送信されたファイルを受信すると (ステップ S213)、クライアント 200 に受信したファイルを送信する (ステップ S214)。クライアント 200 が送信されたファイルを受信すると (ステップ S215)、読み出し処理は完了する。

【0153】

ここで、本発明の実施の形態のデータ読み出し処理についてさらに具体的に説明する。一例として、データ保管サーバ 300C がクライアント 200 からファイル "0012.dat" の読み出し要求を受信したとする。このとき、データ保管サーバ 300C の演算処理装置 380 は、ファイル "0012.dat" のデータ配置管理情報格納サーバを特定する。特定されたデータ配置管理情報格納サーバをデータ保管サーバ 300A とすると、データ保管サーバ 300C はデータ保管サーバ 300A にデータ配置情報の送信を要求する。

20

【0154】

データ保管サーバ A の演算処理装置 380 は、データ配置管理情報 3006 を参照し、ファイル "0012.dat" の格納場所を特定する (図 14)。その結果、ファイル "0012.dat" は、"192.168.0.2" で識別されるデータ保管サーバ 300B の記憶領域 "/mount/data3" に格納されていることがわかる。

30

【0155】

データ保管サーバ 300C の演算処理装置 380 は、ファイルの保管場所として特定されたデータ保管サーバ 300B にファイル "0012.dat" の読み出し及び送信を要求し、データ保管サーバ 300B によって読み出したファイルを受信する。

【0156】

図 23 は、本発明の実施の形態の管理計算機 500 に格納された構成情報を更新する手順を示すフローチャートである。

【0157】

管理計算機 500 の演算処理装置 580 は、構成情報更新サービスプログラム 5004 を実行することによって、定期的に、ストレージ装置 100 及びデータ保管サーバ 300 に構成情報送信要求を送信する (ステップ S301)。

40

【0158】

ストレージ装置 100 のストレージコントローラ 180 又はデータ保管サーバ 300 の演算処理装置 380 は、管理計算機 500 から構成情報送信要求を受信すると、要求された更新情報を管理計算機 500 に送信する (ステップ S302)。

【0159】

管理計算機 500 の演算処理装置 580 は、ストレージ装置 100 及びデータ保管サーバ 300 から送信された構成情報に基づいて、ストレージ装置論理記憶ユニット構成情報 5001 又はデータ保管サーバデータ保管記憶領域構成情報 5002 を最新の状態に更新

50

する（ステップ S 3 0 3）。

【 0 1 6 0 】

図 2 4 は、本発明の実施の形態の管理計算機 5 0 0 のストレージ装置論理記憶ユニット稼働情報 5 0 0 5 を更新するための手順を示すフローチャートである。本処理は、図 1 8 に示した処理の概要において、稼働状況監視処理（ステップ S 1）に対応する。

【 0 1 6 1 】

管理計算機 5 0 0 の演算処理装置 5 8 0 は、稼働情報更新サービスプログラム 5 0 0 6 を実行することによって、定期的に、ストレージ装置 1 0 0 に稼働情報送信要求を送信する（ステップ S 4 0 1）。

【 0 1 6 2 】

ストレージ装置 1 0 0 のストレージコントローラ 1 8 0 は、管理計算機 5 0 0 から稼働情報送信要求を受信すると、記憶領域稼働情報 1 0 0 6 に記録された稼働情報を、論理記憶ユニット構成情報 1 0 0 4 に基づいて、論理記憶ユニット単位に集計する（ステップ S 4 0 2）。さらに、論理記憶ユニット単位に集計された稼働情報を管理計算機 5 0 0 に送信する（ステップ S 4 0 3）。

【 0 1 6 3 】

管理計算機 5 0 0 の演算処理装置 5 8 0 は、ストレージ装置 1 0 0 から送信された稼働情報に基づいて、ストレージ装置論理記憶ユニット稼働情報 5 0 0 5 に更新する（ステップ S 4 0 4）。

【 0 1 6 4 】

図 2 5 から図 3 1 までに示した図は、本発明の実施の形態のデータ保管システムの構成を変更する手順を示すフローチャートである。また、本処理の具体例を図 8 を参照しながら適宜説明する

図 2 5 に示す手順は、具体的に変更する構成の内容を決定し、データ配置管理情報 3 0 0 6 を複製する手順である。

【 0 1 6 5 】

管理計算機 5 0 0 の演算処理装置 5 8 0 は、データ保管記憶領域配置計算プログラム 5 0 0 9 を実行することによって、データ保管記憶領域（論理記憶ユニット 1 1 0 0）の配置先であるデータ保管サーバ 3 0 0 を選定し、変更後の構成を決定する（ステップ S 5 0 1）。

【 0 1 6 6 】

具体的には、ストレージ装置論理記憶ユニット稼働情報 5 0 0 5 を参照し、論理記憶ユニット 1 1 0 0 の負荷状態に基づいて、データ保管記憶領域に順位づけを行う。次に、データ保管サーバ性能情報 5 0 0 7 を参照し、データ入出力性能ランク 5 0 0 7 1 の高いデータ保管サーバ 3 0 0 に高負荷の論理記憶ユニット 1 1 0 0 を割り当てる。同様に、データ入出力性能ランク 5 0 0 7 1 の低いデータ保管サーバ 3 0 0 に低負荷の論理記憶ユニット 1 1 0 0 を割り当てるように変更後の構成を決定する。

【 0 1 6 7 】

さらに具体的に説明すると、図 1 2 に示した記憶領域稼働情報及び図 1 6 に示したストレージ装置論理記憶ユニット稼働情報 5 0 0 5 を参照すると、“LD - 0 2”によって構成される論理記憶ユニット 1 1 0 0 B の負荷が高く、“LD - 0 4”によって構成される論理記憶ユニット 1 1 0 0 D の負荷が低い。

【 0 1 6 8 】

一方、図 1 7 に示したデータ保管サーバ性能情報 5 0 0 7 を参照すると、データ入出力性能ランク 5 0 0 7 1 について、“1 9 2 . 1 6 8 . 0 . 3”で識別されるデータ保管サーバ 3 0 0 C が、“1 9 2 . 1 6 8 . 0 . 1”のデータ保管サーバ 3 0 0 A よりも順位が高くなっている。

【 0 1 6 9 】

したがって、データ保管サーバ“1 9 2 . 1 6 8 . 0 . 3”に論理記憶ユニット 1 1 0 0 A を、データ保管サーバ“1 9 2 . 1 6 8 . 0 . 1”に論理記憶ユニット 1 1 0 0 D を

10

20

30

40

50

配置することによって性能と負荷とのバランスが改善されると判断することができる。

【0170】

続いて、管理計算機500の演算処理装置580は、すべてのデータ保管サーバ300に対してステップS503からステップS508までの処理を繰り返す(ステップS502)。

【0171】

管理計算機500の演算処理装置580は、データ配置管理情報更新プログラム5010を実行することによって、各データ保管サーバ300にデータ配置管理情報3006の更新を要求する(ステップS503)。このとき、各データ保管サーバ300に送信される更新要求には、ステップS501の処理で選定されたデータ保管記憶領域の配置情報が含まれる。

10

【0172】

データ保管サーバ300の演算処理装置380は、データ配置管理情報3006の更新要求を受信すると(ステップS504)、まず、保持しているデータ配置管理情報3006を複製する(ステップS505)。

【0173】

データ保管サーバ300の演算処理装置380は、データ配置管理プログラム3004によって、ステップS504の処理で受信した更新要求に基づいて、ステップS505の処理で複製されたデータ配置管理情報3006を更新する(ステップS506)。

【0174】

具体的には、図14に示したデータ配置管理情報3006のデータ保管サーバ識別情報30062及びデータ保管記憶領域識別情報30063を、ステップS501の処理において前述した内容で更新する。例えば、ファイル“0011.dat”のエントリをデータ保管サーバ“192.168.0.3”、データ保管記憶領域識別情報“/mount/data4”に更新する。その結果、データ保管サーバ300は、後述する図29のステップS545の処理以降、“0011.dat”の配置場所を“/mount/data4”と解釈する。

20

【0175】

データ保管サーバ300は、S506の処理が終了すると、管理計算機500に完了通知を送信する(ステップS507)。その後、管理計算機500は、データ保管サーバ300から完了通知を受信する(ステップS508)。

30

【0176】

図26に示す手順は、変更後の構成に基づいて、データ保管サーバ300によるストレージ装置100へのアクセスが許可されるように設定する手順である。本処理は、図18に示した処理の概要において、アクセス許可処理(ステップS2)に対応する。

【0177】

管理計算機500の演算処理装置580は、ステップS501の処理で選定されたデータ保管記憶領域3100に対応する論理記憶ユニット1100を提供するストレージ装置100に対し、ステップS512からステップS516までの処理を実行する(ステップS511)。

40

【0178】

管理計算機500の演算処理装置580は、データ配置管理プログラム5003によって、ステップS501の処理で選定された変更後のデータ保管サーバ300による当該論理記憶ユニット1100へのアクセスが許可されるように、ストレージ装置100に論理ユニットアクセス許可要求を送信する(ステップS512)。

【0179】

ストレージ装置100のストレージコントローラ180は、論理ユニットアクセス許可要求を受信すると(ステップS513)、変更後のデータ保管サーバ300による論理記憶ユニットへのアクセスが許可されるように論理記憶ユニット構成情報1004を更新する(ステップS514)。

50

【0180】

具体的に説明すると、図11に示した論理記憶ユニット構成情報1004のうち、論理記憶領域“LD-02”のエントリのアクセス許可通信インタフェース識別情報10044に、“50:00:01:1E:0A:E8:A4”を追加する。なお、追加された“50:00:01:1E:0A:E8:A4”は、変更後のデータ保管サーバ300(192.168.0.3)のデータ入出力用通信インタフェース340の識別情報である。

【0181】

ストレージ装置100のストレージコントローラ180は、論理記憶ユニット構成情報1004の更新が完了すると、管理計算機500に完了通知を送信する(ステップS515)。その後、管理計算機500は、ストレージ装置100から完了通知を受信する(ステップS516)。

10

【0182】

図26までに示した処理が完了すると、データ保管処理の受け付けを停止させ、実際に構成を変更する。図27に示す手順は、データ保管サーバ300にデータ保管処理の受け付けを一時的に停止させる手順である。

【0183】

管理計算機500の演算処理装置580は、すべてのデータ保管サーバ300に対し、ステップS532からステップS535までの処理を実行する(ステップS531)。

【0184】

管理計算機500の演算処理装置580は、まず、データ配置管理プログラム5003によって、データ保管処理の受付を一時的に停止するようにデータ保管サーバ300に指示する(ステップS532)。

20

【0185】

データ保管サーバ300の演算処理装置380は、管理計算機500からデータ保管処理の受付停止要求を受信すると、一時的にデータ保管処理の受付を停止する(ステップS533)。

【0186】

その後、データ保管サーバ300の演算処理装置380は、管理計算機500に完了通知を送信する(ステップS534)。そして、管理計算機500は、データ保管サーバ300によって送信された完了通知を受信する(ステップS535)。

30

【0187】

データ保管サーバ300におけるデータ保管処理の受け付けが停止されると、実際に構成が変更される。図28に示す手順は、変更後の構成に基づいて、データ保管サーバ300に論理記憶ユニットをマウントさせる手順である。本処理は、図18に示した処理の概要において、マウント先変更処理(ステップS3)に対応する。

【0188】

管理計算機500の演算処理装置580は、ステップS501の処理で選定されたデータ保管サーバ300に対し、ステップS522からステップS528までの処理を実行する(ステップS521)。

【0189】

管理計算機500の演算処理装置580は、データ配置管理プログラム5003によって、ステップS501の処理で選定された変更後のデータ保管記憶領域3100のマウントを変更後のデータ保管サーバ300に要求する(ステップS522)。

40

【0190】

データ保管サーバ300の演算処理装置380は、管理計算機500からマウント要求を受信すると(ステップS523)、データ保管記憶領域構成情報3003を更新する(ステップS524)。具体的には、図13に示したデータ保管記憶領域構成情報3003において、データ保管記憶領域“/mount/data2”のエントリを、通信インタフェース“50:00:01:1E:0A:E8:04”の論理記憶ユニット“LU-31”に更新する(ステップS524)。

50

【0191】

続いて、データ保管サーバ300の演算処理装置380は、要求されたとおりにマウント処理を実行する(ステップS526)。ステップS526の処理が完了すると、以降、データ保管記憶領域“/mount/data2”の接続先は論理記憶ユニット1100Bから、論理記憶ユニット1100Dに変更される。さらに、容量使用率30034及び転送データ量30035もマウント先の変更にあわせて更新する。

【0192】

その後、データ保管サーバ300の演算処理装置380は、管理計算機500に完了通知を送信する(ステップS527)。そして、管理計算機500は、データ保管サーバ300によって送信された完了通知を受信する(ステップS528)。

10

【0193】

以上の処理によって、変更後の構成に基づいて、データ保管サーバ300に論理記憶ユニットがマウントされる。その後、データ保管処理を再開するためには、データ配置管理情報3006を更新する必要がある。図29に示す手順は、すべてのデータ保管サーバ300のデータ配置管理情報3006に構成の変更を反映させる手順である。本処理は、図18に示した処理の概要において、構成情報更新処理(ステップS4)に対応する。

【0194】

管理計算機500の演算処理装置580は、すべてのデータ保管サーバ300に対し、ステップS542からステップS547までの処理を実行する(ステップS541)。

【0195】

管理計算機500の演算処理装置580は、データ配置管理情報更新プログラム5010によって、データ配置管理情報3006の更新要求をデータ保管サーバ300に送信する(ステップS542)。

20

【0196】

データ保管サーバ300の演算処理装置380は、データ配置管理情報3006の更新要求を受信すると(ステップS543)、ステップS506の処理で作成(複製)された変更後データ配置管理情報を最新の状態に更新する(ステップS544)。この時点では、複製されたデータ配置管理情報3006にのみ変更後の構成が反映されているためである。さらに、データ配置管理情報3006をステップS544の処理で更新された変更後データ配置管理情報に置き換える(ステップS545)。

30

【0197】

なお、ステップS505の処理で一旦複製データを作成し、ステップS506及びステップS544の処理で更新し、ステップS545の処理で置き換える手順は、処理の受け付けを停止する時間を少なくするために実行される。S545の処理でデータ配置管理情報3006を更新する場合、ステップS506及びステップS544の処理を実行する必要はなくなるが、整合性を保つために図25の処理を開始するタイミングでデータ保管処理の受け付けを停止する必要があり、本発明の実施の形態に示した方法と比較して処理の受け付けを停止する時間が長くなる可能性がある。

【0198】

ステップS545の処理が完了すると、データ保管サーバ300の演算処理装置380は、管理計算機500に完了通知を送信する(ステップS546)。そして、管理計算機500は、データ保管サーバ300によって送信された完了通知を受信する(ステップS547)。

40

【0199】

続いて、停止されていたデータ保管処理の受け付けを再開する。図30に示す手順は、データ保管サーバ300にデータ保管処理の受け付けを再開させる手順である。

【0200】

管理計算機500の演算処理装置580は、すべてのデータ保管サーバ300に対し、ステップS552からステップS555までの処理を実行する(ステップS551)。

【0201】

50

管理計算機 500 の演算処理装置 580 は、まず、データ配置管理プログラム 5003 によって、データ保管処理の受付を再開するようにデータ保管サーバ 300 に指示する（ステップ S552）。

【0202】

データ保管サーバ 300 の演算処理装置 380 は、管理計算機 500 からデータ保管処理の受付再開要求を受信すると、データ保管処理の受付を再開する（ステップ S553）。

【0203】

その後、データ保管サーバ 300 の演算処理装置 380 は、管理計算機 500 に完了通知を送信する（ステップ S554）。そして、管理計算機 500 は、データ保管サーバ 300 によって送信された完了通知を受信する（ステップ S555）。

10

【0204】

最後に、構成変更前にデータ保管サーバ 300 からストレージ装置 100 に対するアクセスが許可されていた設定を変更する。図 31 に示す手順は、変更後の構成に基づいて、データ保管サーバ 300 によるストレージ装置 100 へのアクセスが拒否されるように設定する手順である。

【0205】

管理計算機 500 の演算処理装置 580 は、ステップ S501 の処理で変更対象となった論理記憶ユニット 1100 を提供するストレージ装置 100 に対し、ステップ S572 からステップ S576 までの処理を実行する（ステップ S571）。

20

【0206】

管理計算機 500 の演算処理装置 580 は、データ配置管理プログラム 5003 によって、構成変更前に論理記憶ユニット 1100 のアクセスが許可されていたデータ保管サーバ 300 を指定し、当該論理記憶ユニット 1100 へのアクセスが拒否されるようにストレージ装置 100 に指示する（ステップ S572）。

【0207】

ストレージ装置 100 のストレージコントローラ 180 は、論理ユニットアクセス拒否要求を受信すると（ステップ S573）、論理ユニットアクセス拒否要求に基づいて、指定されたデータ保管サーバ 300 による論理記憶ユニットへのアクセスが拒否されるように論理記憶ユニット構成情報 1004 を更新する（ステップ S574）。

30

【0208】

具体的に説明すると、図 11 に示した論理記憶ユニット構成情報 1004 のうち、論理記憶ユニット “LU-12” のエントリのアクセス許可通信インタフェース識別情報 10044 に残存している “50:00:01:1E:0A:E8:A1” を削除する。なお、削除された “50:00:01:1E:0A:E8:A1” は、構成変更前にアクセスが許可されていたデータ保管サーバ 300 (192.168.0.1) のデータ入出力用通信インタフェース 340 の識別情報である。

【0209】

ストレージ装置 100 のストレージコントローラ 180 は、論理記憶ユニット構成情報 1004 の更新が完了すると、管理計算機 500 に完了通知を送信する（ステップ S575）。その後、管理計算機 500 は、ストレージ装置 100 から完了通知を受信する（ステップ S576）。

40

【0210】

以上の処理によって、データ保管システムの構成を変更することができる。最後に、管理者が手動で構成を変更する場合に変更指示を入力する画面について説明する。

【0211】

図 32 は、本発明の実施の形態のデータ保管システムの構成を変更する指示を入力する画面の一例を示す図である。

【0212】

図 25 に示したフローチャートのステップ S501 の処理では、最も負荷の高いデータ

50

保管領域に対して最も性能の高いデータ保管サーバ300を割り当てる例を示したが、管理者によって手動で構成を変更するようにしてもよい。

【0213】

画面上部には、ストレージ装置論理記憶ユニット稼働情報5005に基づいて出力されたデータ保管領域ごとの負荷情報が表示される。画面下部には、割り当てを入れ替えるデータ保管領域を選択するインターフェースが配置されている。

【0214】

管理者は、画面上部に表示された稼働データを閲覧しながらデータ保管記憶領域の最適配置を検討し、入れ替える2つのデータ保管領域を選択する。

【0215】

本発明の実施の形態によれば、長期にわたる運用によってデータ保管システムの各構成に性能差が生じ、データ保管サーバ300の性能とストレージ装置100の負荷状態とのバランスが崩れることによって性能ボトルネックが発生することを回避し、システム全体の性能が低下することを防ぐことができる。

【0216】

また、本発明の実施の形態によれば、データ保管システムをメタデータの更新によって構成を変更するため、ストレージ装置100に格納されたデータの移動を伴わずにシステムを再構成することができる。したがって、システムの構成変更時の運用負荷を低減させ、システムの停止時間を短時間にする事ができる。

【図面の簡単な説明】

【0217】

【図1】本発明の実施の形態のストレージエリアネットワークの構成を表す図である。

【図2】本発明の実施の形態のストレージ装置の構成を示す図である。

【図3】本発明の実施の形態のデータ保管サーバの構成を示す図である。

【図4】本発明の実施の形態の管理計算機の構成を示す図である。

【図5】本発明の実施の形態のストレージ装置のプログラムメモリに格納される制御プログラム及び制御情報の一例を示す図である。

【図6】本発明の実施の形態のデータ保管サーバのプログラムメモリに格納される制御プログラム及び制御情報の一例を示す図である。

【図7】本発明の実施の形態の管理計算機のプログラムメモリに格納される制御プログラム及び制御情報の一例を示す図である。

【図8】本発明の実施の形態のデータ保管システムの構成の一例を示す図である。

【図9】本発明の実施の形態のストレージ装置に格納されるRAIDグループ構成情報の一例を示す図である。

【図10】本発明の実施の形態のストレージ装置に格納される記憶領域構成情報の一例を示す図である。

【図11】本発明の実施の形態のストレージ装置に格納される論理記憶ユニット構成情報の一例を示す図である。

【図12】本発明の実施の形態のストレージ装置に格納される記憶領域稼働情報の一例を示す図である。

【図13】本発明の実施の形態のデータ保管サーバに格納されるデータ保管記憶領域構成情報の一例を示す図である。

【図14】本発明の実施の形態のデータ保管サーバに格納されるデータ配置管理情報の一例を示す図である。

【図15】本発明の実施の形態のデータ保管サーバに格納されるハッシュ計算情報の一例を示す図である。

【図16】本発明の実施の形態の管理計算機に格納されるストレージ装置論理記憶ユニット稼働情報の一例を示す図である。

【図17】本発明の実施の形態の管理計算機に格納されるデータ保管サーバ性能情報の一例を示す図である。

10

20

30

40

50

【図 18】本発明の実施の形態のデータ保管システムの性能を管理するための手順の概要を示す図である。

【図 19】本発明の実施の形態のデータ保管システムによるデータ保管処理の手順を示し、特に、データの保管場所を特定し、データ配置管理情報を更新する手順を示すフローチャートである。

【図 20】本発明の実施の形態のデータ保管システムによるデータ保管処理の手順を示し、特に、実際にファイルを保管する手順を示すフローチャートである。

【図 21】本発明の実施の形態のデータ保管システムに保管されたデータの読み出し処理の手順を示し、特に、データの保管場所を特定する手順を示すフローチャートである。

【図 22】本発明の実施の形態のデータ保管システムに保管されたデータの読み出し処理の手順を示し、特に、実際にファイルを読み出す手順を示すフローチャートである。

【図 23】本発明の実施の形態の管理計算機に格納された構成情報を更新する手順を示すフローチャートである。

【図 24】本発明の実施の形態の管理計算機のストレージ装置論理記憶ユニット稼働情報を更新するための手順を示すフローチャートである。

【図 25】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、具体的に変更する構成の内容を決定し、データ配置管理情報を複製する手順を示すフローチャートである。

【図 26】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、変更後の構成に基づいて、データ保管サーバによるストレージ装置へのアクセスが許可されるように設定する手順を示すフローチャートである。

【図 27】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、データ保管サーバにデータ保管処理の受け付けを一時的に停止させる手順を示すフローチャートである。

【図 28】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、変更後の構成に基づいて、データ保管サーバに論理記憶ユニットをマウントさせる手順を示すフローチャートである。

【図 29】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、すべてのデータ保管サーバのデータ配置管理情報に構成の変更を反映させる手順を示すフローチャートである。

【図 30】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、データ保管サーバにデータ保管処理の受け付けを再開させる手順を示すフローチャートである。

【図 31】本発明の実施の形態のデータ保管システムの構成を変更する手順を示し、特に、変更後の構成に基づいて、データ保管サーバによるストレージ装置へのアクセスが拒否されるように設定する手順を示すフローチャートである。

【図 32】本発明の実施の形態のデータ保管システムの構成を変更する指示を入力する画面の一例を示す図である。

【符号の説明】

【0218】

- 100 ストレージ装置
- 120 磁気記憶装置
- 140 データ入出力用通信インタフェース
- 150 管理用通信インタフェース
- 160 データ入出力用キャッシュメモリ
- 180 ストレージコントローラ
- 190 ストレージコントローラ
- 200 クライアント
- 300 データ保管サーバ
- 320 磁気記憶装置

10

20

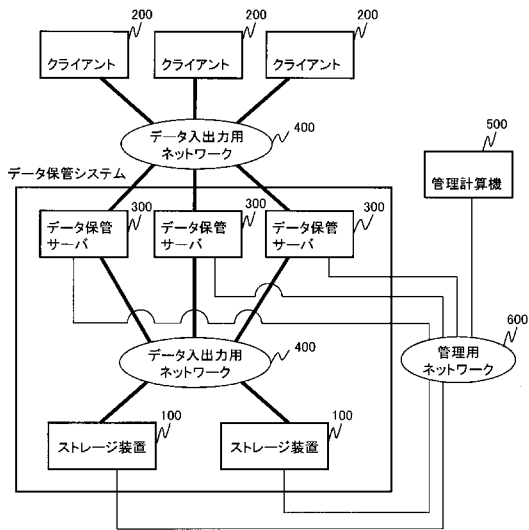
30

40

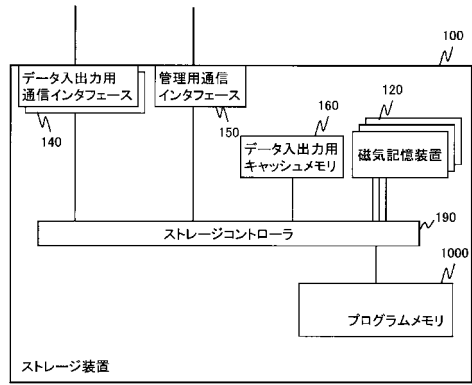
50

3 4 0	データ入出力用通信インタフェース	
3 5 0	管理用通信インタフェース	
3 6 0	データ入出力用キャッシュメモリ	
3 7 0	入力用インタフェース	
3 7 5	出力用インタフェース	
3 8 0	演算処理装置	
3 9 0	通信バス	
4 0 0	データ入出力用ネットワーク	
5 0 0	管理計算機	
5 2 0	磁気記憶装置	10
5 5 0	管理用通信インタフェース	
5 6 0	データ入出力用キャッシュメモリ	
5 7 0	入力用インタフェース	
5 7 5	出力用インタフェース	
5 8 0	演算処理装置	
5 9 0	通信バス	
6 0 0	管理用ネットワーク	
1 0 0 0	プログラムメモリ	
1 0 0 1	記憶領域構成管理プログラム	
1 0 0 2	R A I Dグループ構成情報	20
1 0 0 3	記憶領域構成情報	
1 0 0 4	論理記憶ユニット構成情報	
1 0 0 5	記憶領域稼働監視プログラム	
1 0 0 6	記憶領域稼働情報	
1 1 0 0	論理記憶ユニット	
3 0 0 0	プログラムメモリ	
3 0 0 1	データ入出力プログラム	
3 0 0 2	記憶領域構成管理プログラム	
3 0 0 3	データ保管記憶領域構成情報	
3 0 0 4	データ配置管理プログラム	30
3 0 0 5	データ配置管理情報格納サーバ計算プログラム	
3 0 0 6	データ配置管理情報	
3 0 0 7	ハッシュ計算情報	
3 1 0 0	データ保管記憶領域	
5 0 0 0	プログラムメモリ	
5 0 0 1	ストレージ装置論理記憶ユニット構成情報	
5 0 0 2	データ保管サーバデータ保管記憶領域構成情報	
5 0 0 3	データ配置管理プログラム	
5 0 0 4	構成情報更新サービスプログラム	
5 0 0 5	ストレージ装置論理記憶ユニット稼働情報	40
5 0 0 6	稼働情報更新サービスプログラム	
5 0 0 7	データ保管サーバ性能情報	
5 0 0 8	データ保管サーバ性能情報更新サービスプログラム	
5 0 0 9	データ保管記憶領域配置計算プログラム	
5 0 1 0	データ配置管理情報更新プログラム	

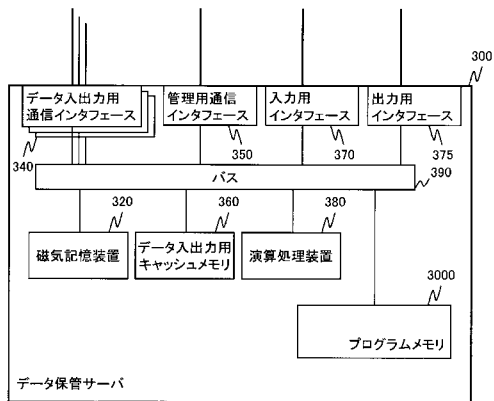
【 図 1 】



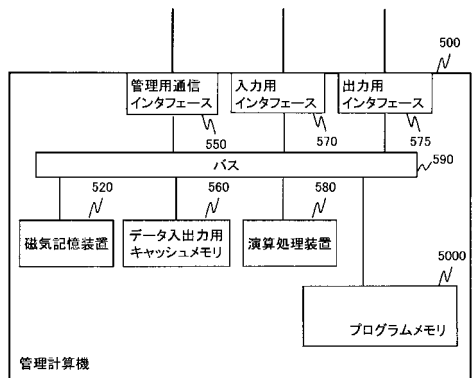
【 図 2 】



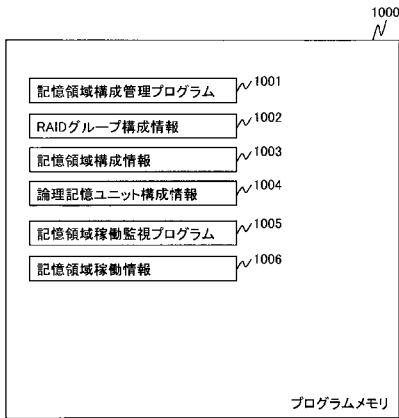
【 図 3 】



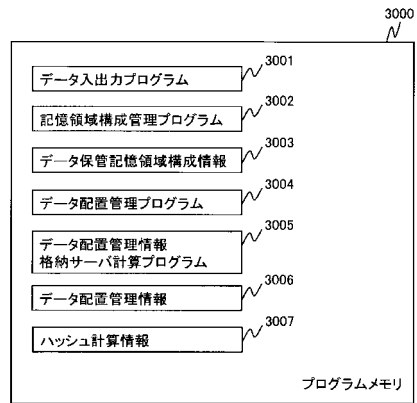
【 図 4 】



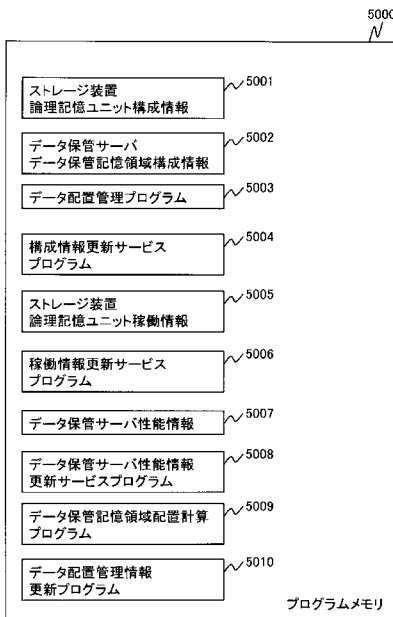
【 図 5 】



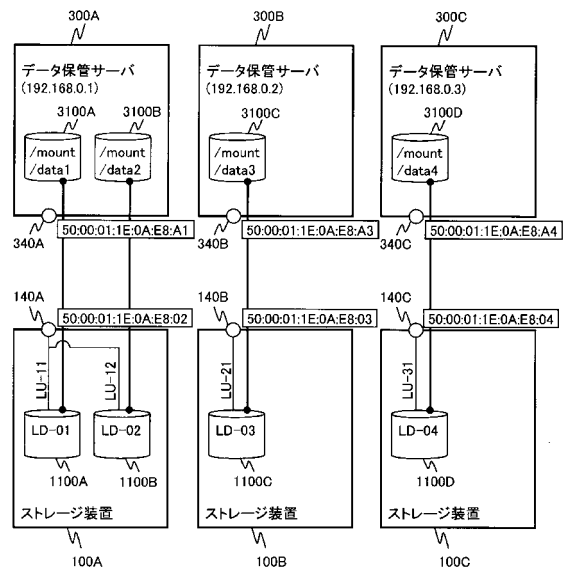
【 図 6 】



【 図 7 】



【 図 8 】



【 図 9 】

RAIDグループ 識別情報	磁気記憶装置識別情報			
	#1	#2	#3	#4
RG-01	HD-01	HD-02	HD-03	HD-04
RG-11	HD-11	HD-12	HD-13	HD-14
...

【 図 1 1 】

通信インタフェース 識別情報	記憶ユニット 識別情報	記憶領域 識別情報	アクセス許可通信 インタフェース識別情報
50:00:01:1E:0A:E8:02	LU-11	LD-01	50:00:01:1E:0A:E8:A1
50:00:01:1E:0A:E8:02	LU-12	LD-02	50:00:01:1E:0A:E8:A1
...

【 図 1 0 】

記憶領域 識別情報	RAID グループ 識別情報	開始 ブロック アドレス	終了 ブロック アドレス
LD-01	RG-01	0x0001	0x0100
LD-02	RG-02	0x0101	0x0200
LD-03	RG-02	0x0201	0x0300
LD-04	RG-03	0x0101	0x0500
...

【 図 1 2 】

Time	LD-01	LD-02	LD-03	...
2008.01.01 0:00	15	2	1	...
2008.01.01 0:01	0	24	2	...
2008.01.01 0:02	23	56	0	...
...

【 図 1 3 】

データ保管記憶領域 識別情報	通信インタフェース 識別情報	記憶ユニット 識別情報	容量使用率	転送データ量
/mount/data1	50:00:01:1E:0A:E8:02	LU-11	40%	10
/mount/data2	50:00:01:1E:0A:E8:02	LU-12	32%	24
...

【 図 1 5 】

サーバ識別情報	サーバ計算情報
192.168.0.1	0
192.168.0.2	1
192.168.0.3	2
...	...

【 図 1 4 】

ファイル識別情報	データ保管サーバ 識別情報	データ保管記憶領域 識別情報
0011.dat	192.168.0.1	/mount/data2
0012.dat	192.168.0.2	/mount/data3
...

【 図 1 6 】

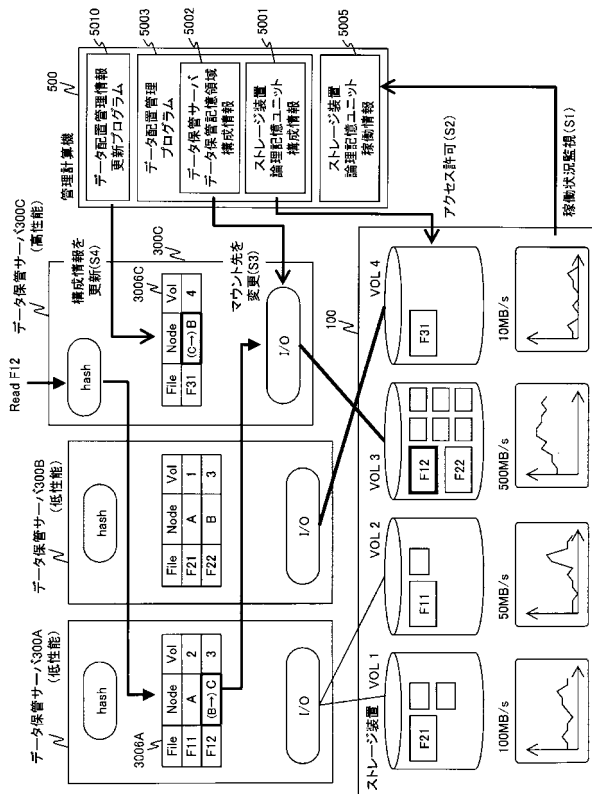
Time	50:00:01:1E:0A:E8:02 LU-11	50:00:01:1E:0A:E8:02 LU-12	...
2008.01.01 0:00	15	2	...
2008.01.01 0:01	0	24	...
2008.01.01 0:02	23	56	...
...

【図 17】

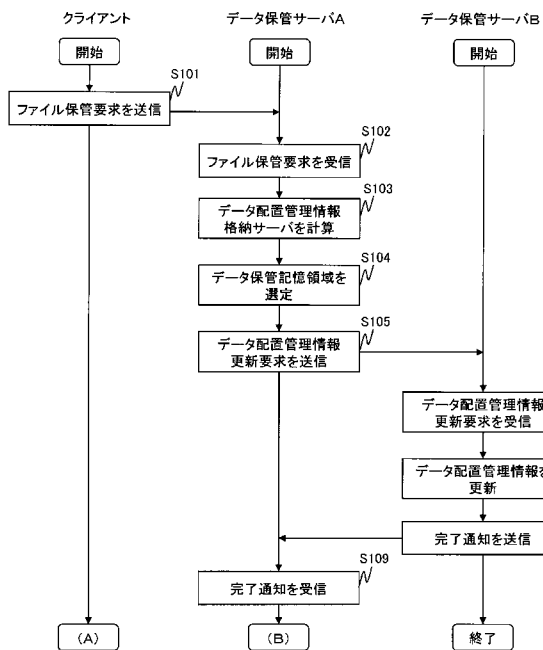
	50071	50072	50073
	データ入出力性能ランク	データ保管サーバ識別情報	通信インタフェース識別情報
1	192.168.0.3		50:00:01:1E:0A:E8:A4
2	192.168.0.1		50:00:01:1E:0A:E8:A1
2	192.168.0.2		50:00:01:1E:0A:E8:A3
...

5007

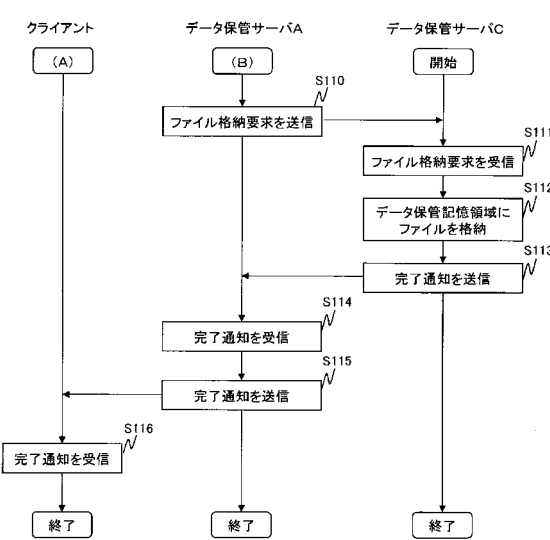
【図 18】



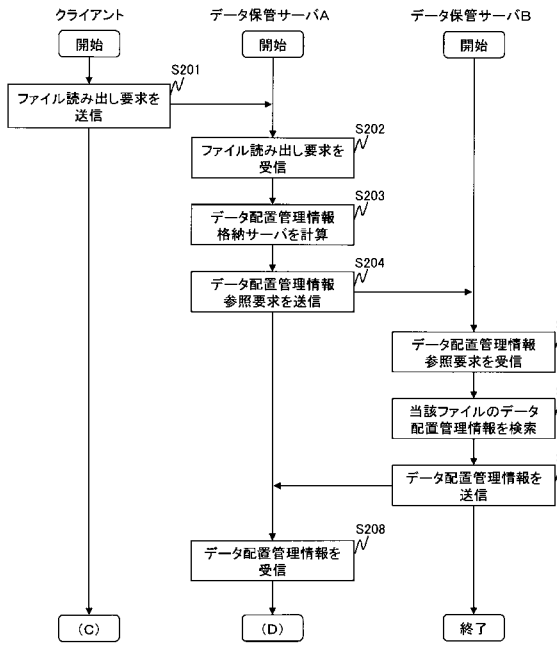
【図 19】



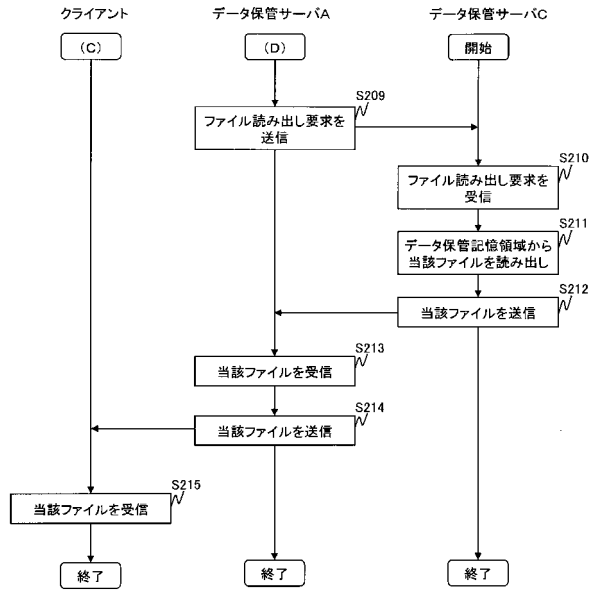
【図 20】



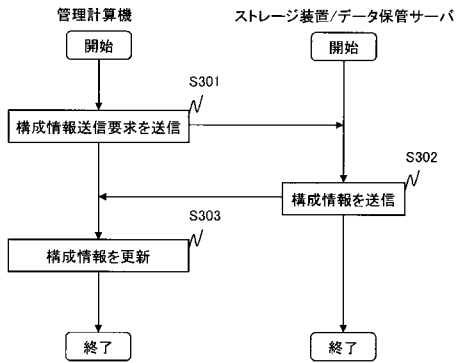
【図 2 1】



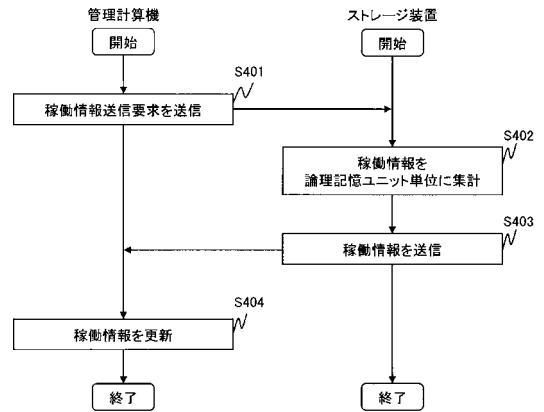
【図 2 2】



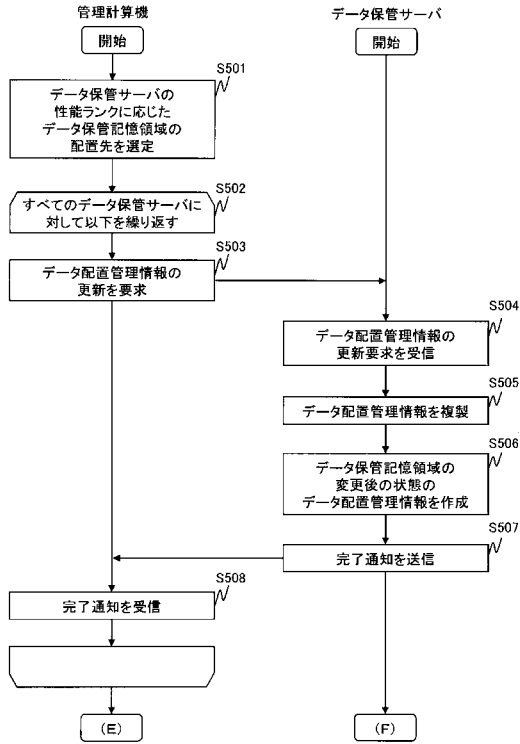
【図 2 3】



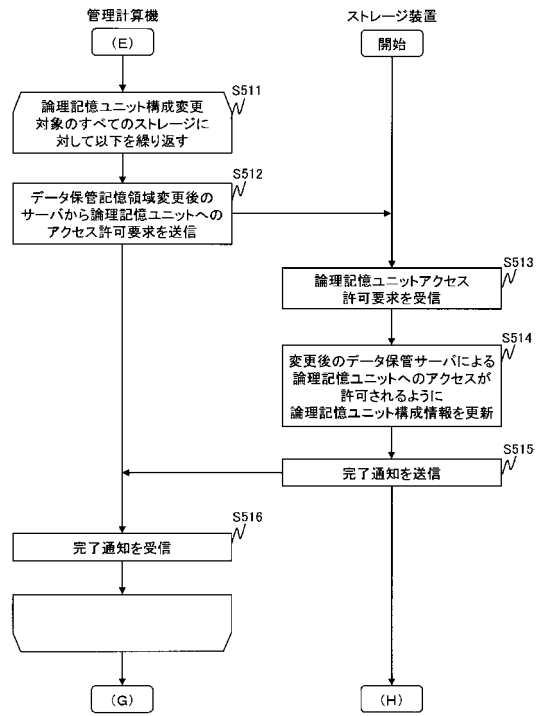
【図 2 4】



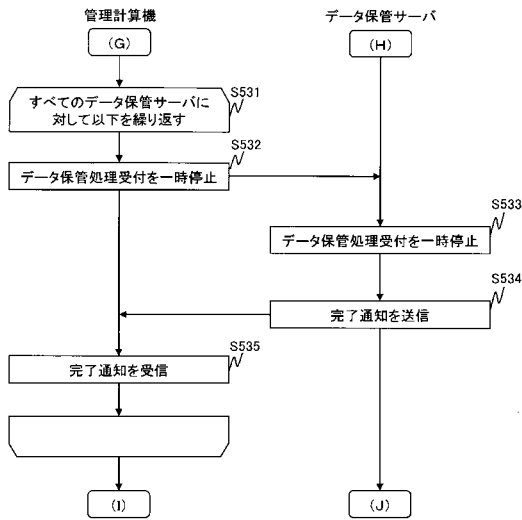
【図 25】



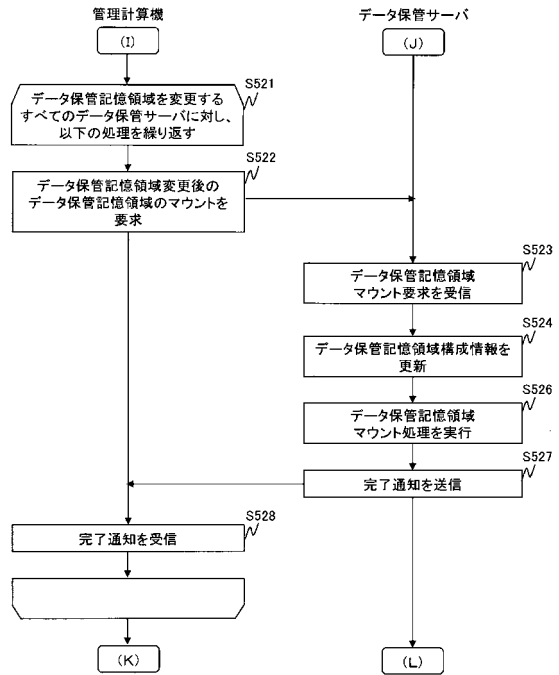
【図 26】



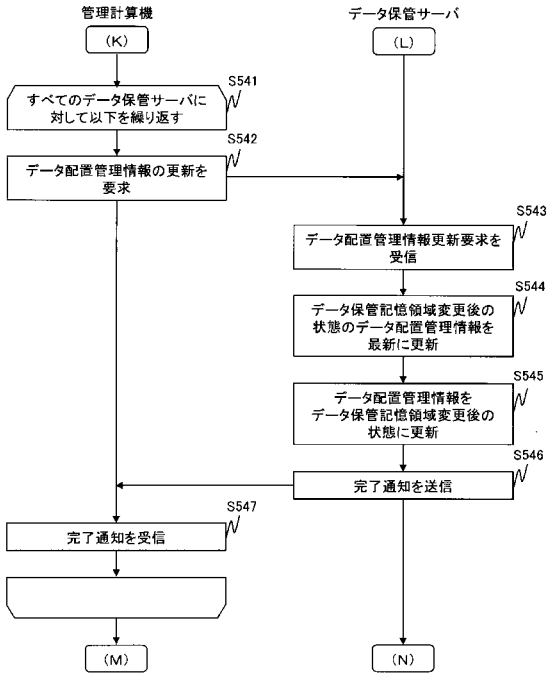
【図 27】



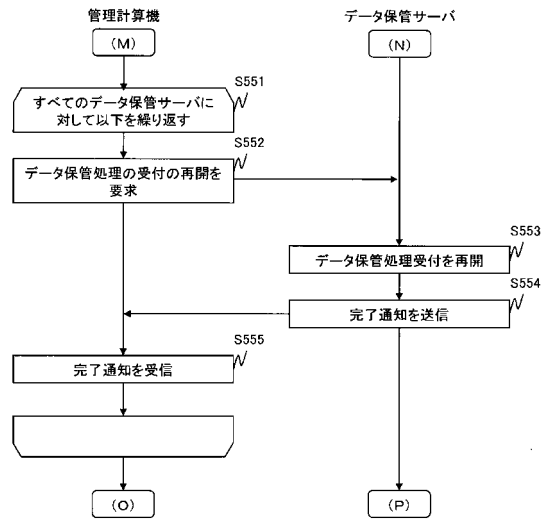
【図 28】



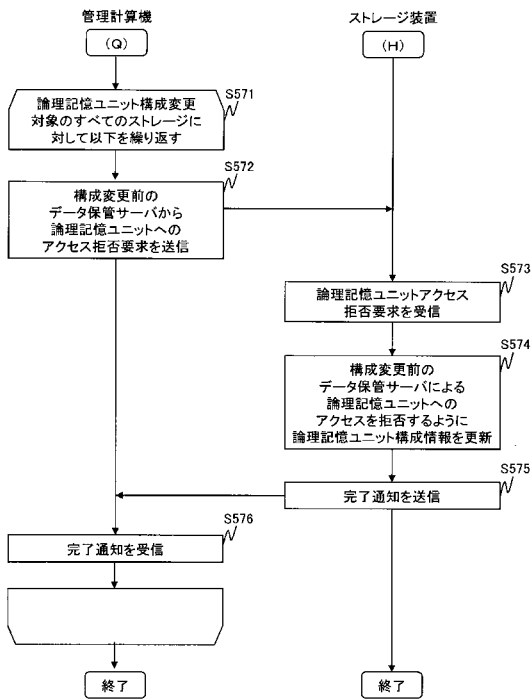
【図 29】



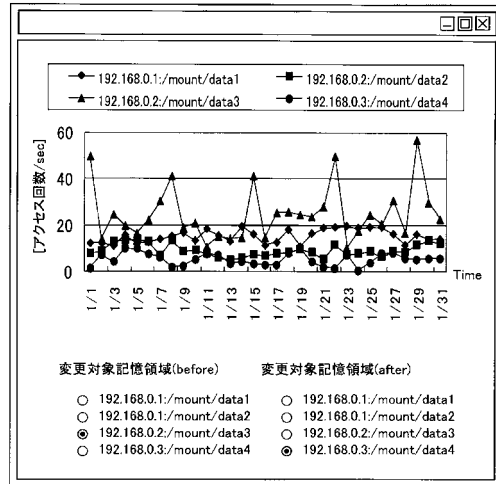
【図 30】



【図 31】



【図 32】



フロントページの続き

(72)発明者 那須 弘志

神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所システム開発研究所内

(72)発明者 平岩 友理

神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所システム開発研究所内

Fターム(参考) 5B082 CA11 CA20