

(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2014-130420

(P2014-130420A)

(43) 公開日 平成26年7月10日(2014.7.10)

(51) Int.Cl.		F I		テーマコード (参考)
G06F 12/00	(2006.01)	G06F 12/00	501B	5B045
G06F 15/17	(2006.01)	G06F 15/17	630A	
G06F 13/10	(2006.01)	G06F 13/10	340A	
G06F 3/08	(2006.01)	G06F 3/08	H	
G06F 3/06	(2006.01)	G06F 3/06	302A	
審査請求 未請求 請求項の数 10 O L (全 21 頁)				

(21) 出願番号 特願2012-286729 (P2012-286729)
 (22) 出願日 平成24年12月28日 (2012.12.28)

(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 100114236
 弁理士 藤井 正弘
 (74) 代理人 100075513
 弁理士 後藤 政喜
 (72) 発明者 杉本 健
 東京都国分寺市東恋ヶ窪一丁目280番地
 株式会社日立製作所中央研究所内
 (72) 発明者 近藤 伸和
 東京都国分寺市東恋ヶ窪一丁目280番地
 株式会社日立製作所中央研究所内
 Fターム(参考) 5B045 BB02 BB12 BB28 BB30

(54) 【発明の名称】 計算機システム及び計算機の制御方法

(57) 【要約】

【課題】サーバに搭載された不揮発性メモリと共有ストレージ間の階層化制御と、サーバ間のデータ配置の最適化制御を動的に行う。

【解決手段】複数のサーバで共有されるデータを格納する共有ストレージと、複数のサーバ及び共有ストレージを管理する管理サーバと、を備え、サーバが不揮発性メモリを備えて前記共有ストレージのデータの一部を格納し、不揮発性メモリに格納されたデータのアクセス状況を格納して第1のアクセス履歴情報を生成し、不揮発性メモリに格納されたデータと共有ストレージに格納されたデータの対応関係を示す格納位置情報に基づいて、不揮発性メモリまたは他のサーバの不揮発性メモリあるいは前記共有ストレージからデータの読み込みまたは書き込みを行う。

【選択図】 図4

図124、224 不揮発性メモリ格納位置情報

共有ストレージセクタ番号	データ格納位置	RDMA利用可否
0 - 419,430,400	サーバ#0 Vol#0 0 - 419,430,400	○
419,430,400 - 838,860,800	サーバ#1 Vol#0 838,860,800 - 1,258,291,200	○
838,860,800 - 1,258,291,200	サーバ#0 Vol#1 0 - 419,430,400	○
その他	共有ストレージ	

【特許請求の範囲】**【請求項 1】**

プロセッサとメモリを備えた複数のサーバと、前記複数のサーバで共有されるデータを格納する共有ストレージと、前記複数のサーバと前記共有ストレージとを接続するネットワークと、前記複数のサーバ及び前記共有ストレージを管理する管理サーバとを備えた計算機システムであって、

前記サーバは、

前記共有ストレージのデータの一部を格納する 1 以上の不揮発性メモリと、

前記ネットワークを介して他のサーバとの間で前記不揮発性メモリのデータを相互に読み書きするインターフェースと、

前記不揮発性メモリに格納されたデータのアクセス状況を格納した第 1 のアクセス履歴情報と、

前記不揮発性メモリに格納されたデータと、前記共有ストレージに格納されたデータの対応関係を保持する格納位置情報と、

前記格納位置情報に基づいて、前記不揮発性メモリまたは前記インターフェースを介して他のサーバの不揮発性メモリあるいは前記共有ストレージからデータの読み込みまたは書き込みを行う第 1 の管理部と、を有し、

前記管理サーバは、

前記第 1 のアクセス履歴情報を前記サーバ毎に取得し、各第 1 のアクセス履歴情報を集約した第 2 のアクセス履歴情報と、

前記第 2 のアクセス履歴情報に基づいて、前記サーバ毎に不揮発性メモリへ配置するデータをそれぞれ決定する第 2 の管理部と、を有することを特徴する計算機システム。

【請求項 2】

請求項 1 に記載の計算機システムであって、

前記第 1 のアクセス履歴情報は、

前記不揮発性メモリに格納されたデータの読み込み回数及び書き込み回数に加えて、前記共有ストレージに格納されたデータの読み込み回数及び書き込み回数を含むことを特徴する計算機システム。

【請求項 3】

請求項 1 に記載の計算機システムであって、

前記第 1 のアクセス履歴情報は、

前記不揮発性メモリに予め設定した履歴の取得単位ごとに前記アクセス状況が格納されることを特徴する計算機システム。

【請求項 4】

請求項 1 に記載の計算機システムであって、

前記第 2 の管理部が、

前記サーバ毎に前記不揮発性メモリへ格納するデータを予め設定し、当該設定に従って前記サーバ毎に前記共有ストレージから不揮発性メモリへ格納するデータを指令することを特徴する計算機システム。

【請求項 5】

請求項 1 に記載の計算機システムであって、

前記第 1 の管理部は、

当該サーバの不揮発性メモリと他のサーバの不揮発性メモリ間でのデータの送受信をリモート DMA を用いて実行し、

前記第 2 の管理部は、

前記第 2 のアクセス履歴情報に基づいて、前記サーバ毎に不揮発性メモリへ配置するデータをそれぞれ決定し、当該配置に従って前記サーバ毎に前記共有ストレージから不揮発性メモリへデータを格納する際には、前記サーバに対して前記リモート DMA を一時的に禁止することを特徴する計算機システム。

【請求項 6】

プロセッサとメモリを備えた複数のサーバと、前記複数のサーバで共有されるデータを格納する共有ストレージと、前記複数のサーバと前記共有ストレージとを接続するネットワークと、前記複数のサーバ及び前記共有ストレージを管理する管理サーバと、を備えてサーバに格納するデータを制御する方法であって、

前記サーバが、1以上の不揮発性メモリを備えて前記共有ストレージのデータの一部を格納する第1のステップと、

前記サーバが、前記不揮発性メモリに格納されたデータのアクセス状況を格納して第1のアクセス履歴情報を生成する第2のステップと、

前記サーバが、前記不揮発性メモリに格納されたデータと、前記共有ストレージに格納されたデータの対応関係を示す格納位置情報を保持し、前記格納位置情報に基づいて、前記不揮発性メモリまたは他のサーバの不揮発性メモリあるいは前記共有ストレージからデータの読み込みまたは書き込みを行う第3のステップと、

前記管理サーバが、前記第1のアクセス履歴情報を前記サーバ毎に取得し、各第1のアクセス履歴情報を集約して第2のアクセス履歴情報を生成する第4のステップと、

前記管理サーバが、前記第2のアクセス履歴情報に基づいて、前記サーバ毎に不揮発性メモリへ配置するデータをそれぞれ決定する第5のステップと、
を含むことを特徴する計算機の制御方法。

【請求項7】

請求項6に記載の計算機の制御方法であって、

前記第1のアクセス履歴情報は、前記不揮発性メモリに格納されたデータの読み込み回数及び書き込み回数に加えて、前記共有ストレージに格納されたデータの読み込み回数及び書き込み回数を含むことを特徴する計算機の制御方法。

【請求項8】

請求項6に記載の計算機の制御方法であって、

前記第2のステップは、

前記不揮発性メモリに予め設定した履歴の取得単位ごとに前記アクセス状況を前記第1のアクセス履歴情報に格納することを特徴する計算機の制御方法。

【請求項9】

請求項6に記載の計算機の制御方法であって、

前記第5のステップは、

前記サーバ毎に前記不揮発性メモリへ格納するデータを予め設定し、当該設定に従って前記サーバ毎に前記共有ストレージから不揮発性メモリへ格納するデータを指令することを特徴する計算機の制御方法。

【請求項10】

請求項6に記載の計算機の制御方法であって、

前記第3のステップは、

当該サーバの不揮発性メモリと他のサーバの不揮発性メモリ間でのデータの送受信をリモートDMAを用いて実行し、

前記第5のステップは、

前記第2のアクセス履歴情報に基づいて、前記サーバ毎に不揮発性メモリへ配置するデータをそれぞれ決定し、当該配置に従って前記サーバ毎に前記共有ストレージから不揮発性メモリへデータを格納する際には、前記サーバに対して前記リモートDMAを一時的に禁止することを特徴する計算機の制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、不揮発性メモリを搭載した複数の計算機で、共有ストレージを利用する技術に関する。

【背景技術】

【0002】

10

20

30

40

50

サーバに搭載する記憶デバイスは、ほぼ100%HDD(Hard Disk Drive)であったが、近年、フラッシュ等の不揮発性メモリのサーバへの搭載が進んでいる。例えば、サーバにPCI Express(以下、PCIe)インターフェースで直結するタイプのフラッシュ(PCIe-SSD)が2011年頃より登場し、徐々に普及している。今後は、MRAM(Magnetic Random Access Memory)、ReRAM(Resistance Random Access Memory)、STTRAM(Spin Transfer Torque Random Access Memory)、PCM(Phase Change Memory)等の不揮発性メモリも様々な形態でサーバに搭載されてゆくと考えられる。

【0003】

これらの不揮発性メモリは、HDDと比べて高速・小容量という特徴がある。その為、HDDと同じようにサーバ直結のストレージとして利用する方法の他に、共有ストレージのキャッシュやティアリング先として利用することによって、サーバ及び共有ストレージ間のI/O性能の向上を図る活用方法が考えられる。これは、共有ストレージとサーバ搭載の不揮発性メモリ間で階層化を行うことで、システム全体でのI/O性能の強化を進める方法である。

【0004】

サーバ直結の不揮発性メモリを、共有ストレージのキャッシュやティアリング先として利用する場合は、サーバ間で不揮発性メモリ上のデータを相互にコピー可能にすることにより、更にI/O性能の向上を図る事が可能である。すなわち、あるサーバが必要とするデータが既に別サーバ上の不揮発性メモリにある場合、そのサーバは別サーバの不揮発性メモリよりデータを取得することによって、共有ストレージへの負荷を減らす事が出来る。似た例として、例えば特開2003-131944号では、共有ストレージのクラスタ環境において、共有ストレージのDRAM間でのデータのコピーを可能にすることによって、I/O性能の向上を図る解決手段が提示されている。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2003-131944号報

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかし、上述した従来の方法には、以下に述べるような課題が存在する。すなわち、サーバ間でのデータのコピーは、共有ストレージからデータを取得する方法の代替手段に過ぎないため、不揮発性メモリ上にはそのサーバ自身が利用するリソースしか配置されず、システム全体での不揮発性メモリの容量の利用効率が上がらない。例えば、サーバに搭載する不揮発性メモリの一例として、PCIe-SSDを考える。通常ベンダーは、PCIe-SSDのラインアップとして数種類のみ用意する。そのため、例えば550Gbyteと1100Gbyteの2つの容量の種類しか存在しない状況が考えられる。このような場合、あるサーバに対して800GbyteのPCIe-SSDの容量が必要な場合は、1100Gbyteの容量を選び、300Gbyteの容量を無駄にすることとなる。

【0007】

本課題に対する解決案として、サーバに搭載した不揮発性メモリをサーバ間で相互にリード・ライト可能にし、複数の不揮発性メモリを、1つの不揮発性メモリに見えるように仮想化する方法が考えられる。この時、仮想化された不揮発性メモリと共有ストレージ間での階層化と、不揮発性メモリのサーバ間でのデータ配置の最適化の両立が課題となる。後者の課題は、あるサーバが利用するデータは、なるべくそのサーバの不揮発性メモリ上においたほうが、効率が上がることに依る。更に、クラスタ構成で稼働するアプリケーションは時間毎に稼働させる種類が変更される場合もあり、この場合では使用するデータなどが変わる事があり、さらに、数ヶ月～数年に一度は、サーバのリプレース若しくは増強

10

20

30

40

50

が行われることの2点を考えると、動的に不揮発性メモリの制御を行う必要がある。

【0008】

そこで本発明は上記で述べたように、サーバに搭載された不揮発性メモリと共有ストレージ間の階層化制御と、サーバ間のデータ配置の最適化制御を動的に行う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】

【0009】

本発明は、プロセッサとメモリを備えた複数のサーバと、前記複数のサーバで共有されるデータを格納する共有ストレージと、前記複数のサーバと前記共有ストレージとを接続するネットワークと、前記複数のサーバ及び前記共有ストレージを管理する管理サーバとを備えた計算機システムであって、前記サーバは、前記共有ストレージのデータの一部を格納する1以上の不揮発性メモリと、前記ネットワークを介して他のサーバとの間で前記不揮発性メモリのデータを相互に読み書きするインターフェースと、前記不揮発性メモリに格納されたデータのアクセス状況を格納した第1のアクセス履歴情報と、前記不揮発性メモリに格納されたデータと、前記共有ストレージに格納されたデータの対応関係を保持する格納位置情報と、前記格納位置情報に基づいて、前記不揮発性メモリまたは前記インターフェースを介して他のサーバの不揮発性メモリあるいは前記共有ストレージからデータの読み込みまたは書き込みを行う第1の管理部と、を有し、前記管理サーバは、前記第1のアクセス履歴情報を前記サーバ毎に取得し、各第1のアクセス履歴情報を集約した第2のアクセス履歴情報と、前記第2のアクセス履歴情報に基づいて、前記サーバ毎に不揮発性メモリへ配置するデータをそれぞれ決定する第2の管理部と、を有する。

10

20

【発明の効果】

【0010】

したがって、本発明は、サーバに搭載された不揮発性メモリの利用効率がシステム全体として改善されると共に、サーバ毎の不揮発性メモリへのデータ配置が最適化される。これにより、計算機システム全体の高性能化及び低コスト化が両立可能となる。

【図面の簡単な説明】

【0011】

【図1】本発明の実施例を示し、計算機システムの一例を示すブロック図である。

【図2】本発明の実施例を示し、自サーバ内不揮発性メモリ利用情報の一例を示す図である。

30

【図3】本発明の実施例を示し、クラスタ内不揮発性メモリ利用情報の一例を示す図である。

【図4】本発明の実施例を示し、不揮発性メモリ格納位置情報の一例を示す図である。

【図5】本発明の実施例を示し、自サーバ内アクセス履歴情報の一例を示す図である。

【図6】本発明の実施例を示し、クラスタ内アクセス履歴情報の一例を示す図である。

【図7】本発明の実施例を示し、リードを行う場合に、サーバで行われる処理の一例を説明するフローチャートである。

【図8】本発明の実施例を示し、ライトを行う場合に、サーバで行われる処理の一例を説明するフローチャートである。

40

【図9】本発明の実施例を示し、マスターサーバで行われるデータ配置変更の全体の処理を説明するフローチャートである。

【図10】本発明の実施例を示し、マスターサーバ及びサーバで行われる不揮発性メモリへ配置するデータの決定処理の一例を示すフローチャートである。

【図11】本発明の実施例を示し、マスターサーバ及びサーバで行われるデータ配置の更新処理で、サーバの不揮発性メモリへのデータ登録処理及び削除処理の一例を示すフローチャートである。

【図12】本発明の実施例を示し、マスターサーバ及びサーバで行われるサーバの不揮発性メモリへのデータ登録処理の一例を示すフローチャートである。

【図13】本発明の実施例を示し、マスターサーバ及びサーバで行われるサーバの不揮発

50

性メモリのデータ削除処理の一例を示すフローチャートである。

【図14】本発明の実施例を示し、マスターサーバ及びサーバで行われる初期化処理の一例を説明するフローチャートである。

【図15】本発明の実施例を示し、電源遮断時にマスターサーバ及びサーバで行われる処理の一例を説明するフローチャートである。

【発明を実施するための形態】

【0012】

以下、本発明の実施の形態を図面に基づいて詳細に説明する。

【0013】

図1は、本発明を適用した計算機システムのブロック図を示す。図1においてサーバと共有ストレージシステムを含む計算機システムは、1台以上のサーバ100-1~100-nと、マスターサーバ(管理サーバ)200と、サーバ間を接続するサーバ間インターコネクト300(またはネットワーク)と、サーバ100-1~100-nと共有ストレージシステム500間を接続する共有ストレージインターコネクト400(またはネットワーク)と、データを格納する共有ストレージシステム500を有する。サーバ100-1~100-nとマスターサーバ200を接続するサーバ間インターコネクト300としては、例えば、Ethernet(登録商標)やInfinibandに準拠した規格を適用することができる。また、共有ストレージインターコネクト400としては、例えば、Fiber channelに準拠した規格を適用することができる。

10

【0014】

20

サーバ100-1は、プロセッサ110-1と、メモリ120-1と、サーバ100-1とサーバ間インターコネクト300を接続するインターフェース130-1と、サーバ100-1と共有ストレージインターコネクト400を接続するインターフェース131-1と、サーバ100-1と不揮発性メモリ140-1間を接続するインターフェース132-1と、不揮発性メモリ140-1を有する。

【0015】

サーバ100-1と不揮発性メモリ140-1の間は、例えば、PCI-SIG(<http://www.pcisig.com/>)が策定したPCI Express(PCIe)の規格に準拠した規格が用いられているものとする。また、不揮発性メモリ140-1は、例えば、フラッシュメモリなどの記憶素子で構成される。

30

【0016】

サーバ100-1の不揮発性メモリ140-1と、サーバ100-nの不揮発性メモリ140-2は、インターフェース131-1、131-2及び共有ストレージインターコネクト400を介して相互に接続され、RDMA(Remote Dynamic Memory Access)を利用して不揮発性メモリ140-1、140-2間でデータを転送することができる。

【0017】

メモリ120-1上には、自サーバ内不揮発性メモリ管理アプリケーション121-1と、自サーバ内不揮発性メモリ利用情報122-1と自サーバ内アクセス履歴情報123-1と、不揮発性メモリ格納位置情報124-1が読み込まれる。自サーバ内不揮発性メモリ管理アプリケーション121-1は、例えば、共有ストレージシステム500に格納されており、プロセッサ110-1が、メモリ120-1にロードして実行する。

40

【0018】

プロセッサ110-1は、各機能部のプログラムに従って動作することによって、所定の機能を実現する機能部として動作する。例えば、プロセッサ110-1は、自サーバ内不揮発性メモリ管理アプリケーション(プログラム)121-1に従って動作することで自サーバ内不揮発性メモリ管理部として機能する。他のプログラムについても同様である。さらに、プロセッサ110-1は、各プログラムが実行する複数の処理のそれぞれを実現する機能部としても動作する。計算機及び計算機システムは、これらの機能部を含む装置及びシステムである。

50

【0019】

サーバ100-1の各機能を実現するプログラム、テーブル等の情報は、共有ストレージシステム500や不揮発性半導体メモリ、ハードディスクドライブ、SSD(Solid State Drive)等の記憶デバイス、または、ICカード、SDカード、DVD等の計算機読み取り可能な非一時的データ記憶媒体に格納することができる。

【0020】

なお、全てのサーバ100-1~100-nは上記で述べたものと同一のハードウェアとソフトウェアを有するので、サーバ100-1と重複する説明は省略する。また、サーバ100-1~100-nの総称を、添え字のない符号100で示す。他の構成要素の総称についても同様であり、添え字のない符号で示す。

10

【0021】

マスターサーバ200は、プロセッサ210と、メモリ220と、サーバ100-1~100-nとサーバ間インターコネクタ300を接続するインターフェース230を有する。また、メモリ220上に、クラスタ内不揮発性メモリ管理アプリケーション221と、クラスタ内不揮発性メモリ利用情報222とクラスタ内アクセス履歴情報223と、不揮発性メモリ格納位置情報224(図4)を有する。

【0022】

プロセッサ210は、各機能部のプログラムに従って動作することによって、所定の機能を実現する機能部として動作する。例えば、プロセッサ110-11は、クラスタ内不揮発性メモリ管理アプリケーション(プログラム)221に従って動作することでクラスタ内不揮発性メモリ管理部として機能する。他のプログラムについても同様である。さらに、プロセッサ210は、各プログラムが実行する複数の処理のそれぞれを実現する機能部としても動作する。計算機及び計算機システムは、これらの機能部を含む装置及びシステムである。

20

【0023】

マスターサーバ200の各機能を実現するプログラム、テーブル等の情報は、共有ストレージシステム500や不揮発性半導体メモリ、ハードディスクドライブ、SSD(Solid State Drive)等の記憶デバイス、または、ICカード、SDカード、DVD等の計算機読み取り可能な非一時的データ記憶媒体に格納することができる。

【0024】

図2は、図1で符号122-1~122-nで示される自サーバ内不揮発性メモリ利用情報を示す。なお、自サーバ内不揮発性メモリ利用情報の総称を符号122で示す。

30

【0025】

自サーバ内不揮発性メモリ利用情報122は、各サーバ100に搭載されている不揮発性メモリデバイスそれぞれに対して識別子としての番号を付与した不揮発性メモリデバイス番号1221と、それぞれの不揮発性メモリ140の容量1222と、それぞれの不揮発性メモリ140で使用中を示す利用容量1223と、利用されている不揮発性メモリ140のセクタ番号を格納する不揮発性メモリセクタ番号1224と、当該不揮発性メモリセクタ番号1224に対応する共有ストレージセクタ番号1225を含む。なお、利用中の不揮発性メモリセクタ番号1224に対応する共有ストレージシステム500のセクタ番号が存在しない場合には、共有ストレージセクタ番号1225に「利用なし」が設定される。

40

【0026】

図3は、符号221で示されるクラスタ内不揮発性メモリ利用情報を示す。クラスタ内不揮発性メモリ利用情報222は、各サーバ100の識別子を格納するサーバ番号2221と、各サーバ100が保持している不揮発性メモリ140の合計容量を格納する不揮発性メモリ合計容量2222と、サーバ100が利用している不揮発性メモリ140の合計利用容量2223と、それぞれの不揮発性メモリ140に保持された共有ストレージシステム500のセクタ番号を格納する不揮発性メモリ内保持共有ストレージセクタ番号2224を含む。

50

【 0 0 2 7 】

図 4 は、図 1 の符号 1 2 4 - 1 ~ 1 2 4 - n 及び 2 2 4 で示される不揮発性メモリ格納位置情報を示す。不揮発性メモリ格納位置情報 1 2 4、2 2 4 は、共有ストレージシステム 5 0 0 のデータの内、サーバ 1 0 0 - 1 ~ 1 0 0 - n の不揮発性メモリ 1 4 0 に格納されているデータの位置（セクタ番号）を、サーバ 1 0 0 の識別子と、不揮発性メモリ 1 4 0 のセクタ番号に対応付けて管理するテーブルである。

【 0 0 2 8 】

不揮発性メモリ格納位置情報 1 2 4、2 2 4 は、共有ストレージシステム 5 0 0 のセクタ番号を格納する共有ストレージセクタ番号 1 2 4 1 と、サーバ 1 0 0 の識別子と各サーバ 1 0 0 に格納された共有ストレージセクタ番号 1 2 4 1 に対応する不揮発性メモリ 1 4 0 のセクタ番号を格納するデータ格納位置 1 2 4 2 及び、当該データへ R D M A を用いてアクセス可能で有るか否かを示す R D M A 利用可否 1 2 4 3 で構成される。データ格納位置は、当該データが存在する不揮発性メモリ 1 4 0 を有するサーバ番号と、不揮発性メモリ番号（V o l # 0 ~ # n）と不揮発性メモリ 1 4 0 のセクタ番号と、当該データが、共有ストレージシステム 5 0 0 上に存在する事示す情報を有する。R D M A を用いてアクセス可能で有るか否かは、不揮発性メモリ 1 4 0 自体が別サーバ 1 0 0 からの R D M A に対応しているか否かによって定まる他、不揮発性メモリ 1 4 0 へのデータ登録処理時に一時的に利用不可に設定される。R D M A 利用可否 1 2 4 3 は、利用可能であれば「 」が設定される。

10

【 0 0 2 9 】

図 5 は、図 1 の符号 1 2 3 - 1 ~ 1 2 3 - n で示される自サーバ内アクセス履歴情報を示す。自サーバ内アクセス履歴情報 1 2 3 は、アクセス履歴取得単位 1 2 3 1 と、不揮発性メモリ 1 4 0 に格納された共有ストレージセクタ番号 1 2 3 2 及び、サーバ 1 0 0 がアクセス履歴取得単位 1 2 3 1 毎に取得したリード回数 1 2 3 3 とライト回数 1 2 3 4 の情報を有する。アクセス履歴取得単位 1 2 3 1 は、各サーバ 1 0 0 がアクセス回数を記録する不揮発性メモリ 1 4 0 の単位を表し、例えば、ボリューム番号やブロック番号等で表される。なお、リード回数 1 2 3 3 とライト回数 1 2 3 4 は、自サーバ 1 0 0 の不揮発性メモリ 1 4 0 へのアクセス回数と、共有ストレージシステム 5 0 0 へのアクセス回数の和で構成することができる。

20

【 0 0 3 0 】

図 6 は、図 1 の符号 2 2 3 で示されるクラスタ内アクセス履歴情報を示す。クラスタ内アクセス履歴情報 2 2 3 は、アクセス履歴取得単位 2 2 3 1 と共有ストレージセクタ番号 2 2 3 2 及び、各サーバ 1 0 0 が、アクセス履歴取得単位 2 2 3 1 毎に取得したリード回数とライト回数の情報 2 2 3 3（2 2 3 3 R、2 2 3 3 W）、2 2 3 4（2 2 3 4 R、2 2 3 4 W）と、全サーバ 1 0 0 - 1 ~ 1 0 0 - n が合計でアクセス履歴取得単位 2 2 3 1 毎に取得したリード回数 2 2 3 5 R とライト回数 2 2 3 5 W の情報を有する。図 6 の例は、サーバ 1 0 0 が 2 台の例を示す。

30

【 0 0 3 1 】

図 7、図 8 に、本発明におけるリード及びライト動作のフローチャートを示す。以下、図 7、図 8 を用いてリード及びライト時の動作を説明する。

40

【 0 0 3 2 】

図 7 は、サーバ 1 0 0 がリードを行う場合に行われる処理の一例を示すフローチャートである。この処理は、サーバ 1 0 0 で稼働する図示しないアプリケーションや O S がデータの読み出し要求を発行したときに自サーバ内不揮発性メモリ管理アプリケーション 1 2 1 で実行される。

【 0 0 3 3 】

まず、ステップ S 1 0 0 において、サーバ 1 0 0 は、リード先の共有ストレージシステム 5 0 0 のセクタ番号と、不揮発性メモリ格納位置情報 1 2 4 より、読み出し対象のデータが不揮発性メモリ 1 4 0 上に存在するか否かを読み出し、また、読み出し対象が不揮発性メモリ 1 4 0 に存在する場合はデータの格納位置情報を読み出す。さらに、サーバ 1 0

50

0 は自サーバ内アクセス履歴情報 1 2 3 のリード回数 1 2 3 3 の対象エントリをインクリメントする。

【 0 0 3 4 】

次にステップ S 1 0 1 において、サーバ 1 0 0 は読み込み対象のデータが自サーバ 1 0 0 内の不揮発性メモリ 1 4 0 に存在するか否かを判定する。読み出し対象のデータが不揮発性メモリ 1 4 0 存在する場合はステップ S 1 0 2 へ進む。ステップ S 1 0 2 では、自サーバ 1 0 0 内の不揮発性メモリ 1 4 0 よりデータを読み出す。データの読み出し位置は不揮発性メモリ格納位置情報 1 2 4 より取得することができる。

【 0 0 3 5 】

一方、ステップ S 1 0 1 において、不揮発性メモリ格納位置情報 1 2 4 を読み込んで、読み出し対象のデータが自サーバ 1 0 0 内の不揮発性メモリ 1 4 0 に存在しなかった場合は、ステップ S 1 0 3 に進み、サーバ 1 0 0 は読み出し対象のデータが他サーバ 1 0 0 - n (以下、データ格納先サーバ 1 0 0 - n) の不揮発性メモリ 1 4 0 に存在するか否かを判定する。

10

【 0 0 3 6 】

読み出し対象のデータが他のサーバ 1 0 0 - n に存在する場合はステップ S 1 0 4 に進み、不揮発性メモリ格納位置情報 1 2 4 の R D M A 利用可否 1 2 4 3 を参照して、読み出し対象のデータ格納先サーバ 1 0 0 - n で R D M A が利用可能か否かを判定する。

【 0 0 3 7 】

R D M A が利用可能な場合はステップ S 1 0 5 に進み、サーバ 1 0 0 は、読み出し対象のデータ格納したサーバ 1 0 0 - n の不揮発性メモリ 1 4 0 - n から R D M A を用いて該当するデータを読み出す。このとき、データ通信経路としてサーバ間インターコネクト 3 0 0 を用いることができる。

20

【 0 0 3 8 】

その結果、ステップ S 1 0 6 において、読み出し対象のデータ格納先サーバ 1 0 0 - n の不揮発性メモリ 1 4 0 - n からデータが読み出され、データを要求した読み出し元のサーバ 1 0 0 はレスポンスを受信する。このような、他サーバ 1 0 0 - n の不揮発性メモリ 1 4 0 から R D M A を用いてデータを読み出す規格としては例えば、S R P (S C S I R D M A P r o t o c o l) と呼ばれる規格がある。

【 0 0 3 9 】

30

一方、ステップ S 1 0 4 において、読み出し対象のデータ格納先サーバ 1 0 0 - n で R D M A が利用不可能であった場合、ステップ S 1 0 7 に進み、読み出し元のサーバ 1 0 0 はデータ格納先のサーバ 1 0 0 - n に対して不揮発性メモリ 1 4 0 からのデータ読み出しを依頼する。

【 0 0 4 0 】

ステップ S 1 0 8 において、データ格納先サーバ 1 0 0 - n のプロセッサ 1 1 0 - n は、要求されたデータを不揮発性メモリ 1 4 0 - n 若しくは共有ストレージシステム 5 0 0 のいずれかから読み出し、読み出し元のサーバ 1 0 0 へレスポンスを返す。

【 0 0 4 1 】

R D M A が利用不可能の場合は、サーバ 1 0 0 のプロセッサ 1 1 0 間でデータの送受を行う。このとき、データの通信経路としてサーバ間インターコネクト 3 0 0 を用いることができる。

40

【 0 0 4 2 】

ステップ S 1 0 3 の判定において、読み出し対象のデータが他サーバ 1 0 0 - n の不揮発性メモリ 1 4 0 にも存在しなかった場合、サーバ 1 0 0 は読み出し対象のデータが計算機システム全体の不揮発性メモリ 1 4 0 上に存在しないと判断し、ステップ S 1 0 9 に進んで共有ストレージシステム 5 0 0 よりデータを読み出す。

【 0 0 4 3 】

以上の処理により、データを読み出すサーバ 1 0 0 では、読み出し要求を受け付けると、自サーバ内不揮発性メモリ管理アプリケーション 1 2 1 によって自サーバ 1 0 0 - 1 の

50

不揮発性メモリ 140 - 1 または他のサーバ 100 - n の不揮発性メモリ 140 - n から優先してデータを読み込むことで、共有ストレージシステム 500 のデータを効率よく利用することができる。

【0044】

図 8 は、サーバ 100 がライトを行う場合に行われる処理の一例を示すフローチャートである。この処理は、サーバ 100 で稼働する図示しないアプリケーションや OS がデータの書き込み要求を発行したときに自サーバ内不揮発性メモリ管理アプリケーション 121 で実行される。

【0045】

まず、ステップ S200 において、サーバ 100 は、ライト先の共有ストレージシステム 500 のセクタ番号を取得してから、不揮発性メモリ格納位置情報 124 を参照して、書き込み対象のデータが不揮発性メモリ 140 上に存在するか否かを判定し、データの格納位置情報を読み出す。さらに、サーバ 100 は、自サーバ内アクセス履歴情報 123 のライト回数 1234 の対象エントリをインクリメントする。

【0046】

次に、ステップ S201 において、サーバ 100 は書き込み対象のデータが自サーバ 100 内の不揮発性メモリ 140 に存在するか否かを判定する。自サーバ 100 内に存在する場合は、ステップ S202 に進んで、自サーバ 100 内の不揮発性メモリ 140 に書き込み対象のデータを書き込む。書き込み位置は不揮発性メモリ格納位置情報 124 より取得することができる。

【0047】

そして、ステップ S203 において、サーバ 100 は、不揮発性メモリ 140 に書き込んだデータを共有ストレージシステム 500 へ書き込む。

【0048】

一方、ステップ S201 の判定において、書き込み対象のデータが自サーバ 100 内の不揮発性メモリ 140 に存在しなかった場合は、ステップ S204 に進んで、書き込み対象のデータが他サーバ 100 - n の不揮発性メモリ 140 - n に存在するか否かを判定する。書き込み対象のデータが他のサーバ 100 - n に存在する場合は、ステップ S205 に進んで、他のサーバ 100 - n で RDMA が利用可能か否かを判定する。他のサーバ 100 - n で RDMA が利用可能な場合はステップ S206 に進んで、サーバ 100 は RDMA を用いてデータの格納先となる他のサーバ 100 - n の不揮発性メモリ 140 へデータを書き込む。

【0049】

次に、ステップ S207 において、データ格納先のサーバ 100 - n の不揮発性メモリ 140 - n へデータが書き込まれると、データ書き込み元のサーバ 100 はデータを書き込んだ他のサーバ 100 - n のインターフェース 131 - n からレスポンスを受け取る。

【0050】

その後、サーバ 100 はステップ S208 で、他のサーバ 100 - n の不揮発性メモリ 140 - n に書き込んだデータを、共有ストレージシステム 500 へ書き込む。

【0051】

一方、ステップ S205 の判定において、データ格納先の他のサーバ 100 - n で RDMA が利用不可能であった場合、ステップ S209 に進む。ステップ S209 では、書き込み元のサーバ 100 はデータ格納先のサーバ 100 - n に対して不揮発性メモリ 140 のデータを不揮発性メモリ 140 - n へ書き込む依頼を実施する。

【0052】

次に、ステップ S210 において、データ格納先のサーバ 100 - n はサーバ 100 から受信したデータを不揮発性メモリ 140 へ書き込み、書き込み元のサーバ 100 へレスポンスを返す。レスポンスを受信した後、書き込み元のサーバ 100 は、ステップ S211 において、他のサーバ 100 - n に依頼して不揮発性メモリ 140 - n に書き込んだデータを、共有ストレージシステム 500 へ書き込む。

10

20

30

40

50

【 0 0 5 3 】

R D M A が利用不可能の場合は、上記読みだしと同様であり、サーバ 1 0 0 のプロセッサ 1 1 0 間でデータの送受を行う。このとき、データの通信経路としてサーバ間インターコネクト 3 0 0 を用いることができる。

【 0 0 5 4 】

一方、上記ステップ S 2 0 4 の判定において、書き込み対象のデータが他サーバ 1 0 0 - n の不揮発性メモリ 1 4 0 にも存在しなかった場合、サーバ 1 0 0 は、書き込み対象のデータが計算機システム全体の不揮発性メモリ 1 4 0 上に存在しないと判断し、ステップ S 2 1 2 において共有ストレージシステム 5 0 0 へデータを書き出す。

【 0 0 5 5 】

以上の処理により、データを書き込むサーバ 1 0 0 では、書き込み要求を受け付けると、自サーバ内不揮発性メモリ管理アプリケーション 1 2 1 によって自サーバ 1 0 0 - 1 の不揮発性メモリ 1 4 0 - 1 または他のサーバ 1 0 0 - n の不揮発性メモリ 1 4 0 - n へ優先してデータを書き込むことで、共有ストレージシステム 5 0 0 のデータを効率よく利用することができる。

【 0 0 5 6 】

図 7、図 8 より、本計算機システムにおけるリード及びライト処理のプロセッサ 1 1 0 のオーバーヘッドは、不揮発性メモリ格納位置情報 1 2 4 より読み出し対象のデータまたは書き込み対象のデータの格納位置を読み出す処理と、自サーバ内アクセス履歴情報 1 2 3 の対象エントリをインクリメントする処理の 2 つのみであり、I / O 処理のオーバーヘッドとして問題にならない程度となる。

【 0 0 5 7 】

次に、図 9、図 1 0、図 1 1、図 1 2、図 1 3 において、マスターサーバ 2 0 0 が各サーバ 1 0 0 の不揮発性メモリ 1 4 0 へデータを登録する処理を示す。

【 0 0 5 8 】

マスターサーバ 2 0 0 は、例えば一定時間毎や予め定められたタイミングで、サーバ 1 0 0 上の不揮発性メモリ 1 4 0 へのデータ格納の配置を変更するフローチャートを実行する。図 9 は、上記所定のタイミングで実行される処理の全体を示したフローチャートである。これらの処理は、マスターサーバ 2 0 0 のクラスタ内不揮発性メモリ管理アプリケーション 2 2 1 と、各サーバ 1 0 0 の自サーバ内不揮発性メモリ管理アプリケーション 1 2 1 によって実行される。

【 0 0 5 9 】

まず、ステップ S 3 0 0 において、マスターサーバ 2 0 0 は各サーバ 1 0 0 の不揮発性メモリ 1 4 0 の新たなデータ配置を決定する。次にステップ S 3 0 1 において、マスターサーバ 2 0 0 は、ステップ S 3 0 0 において決定された新たな配置に従って、各サーバ 1 0 0 に不揮発性メモリ 1 4 0 へのデータ登録または削除を指示する。

【 0 0 6 0 】

図 1 0 は、図 9 のステップ S 3 0 0 で示される、マスターサーバ 2 0 0 が各サーバ 1 0 0 の不揮発性メモリ 1 4 0 毎に新たなデータ配置を決定する処理の詳細なフローチャートを示す。この処理は、マスターサーバ 2 0 0 のクラスタ内不揮発性メモリ管理アプリケーション 2 2 1 と、各サーバ 1 0 0 の自サーバ内不揮発性メモリ管理アプリケーション 1 2 1 によって実行される。なお、以下のフローチャートも同様である。

【 0 0 6 1 】

まず、マスターサーバ 2 0 0 は、ステップ S 4 0 1 において、全サーバ 1 0 0 に対して自サーバ内アクセス履歴情報 1 2 3 の送信を要求する。これに対し各サーバ 1 0 0 は、ステップ S 4 0 2 において、自サーバ内アクセス履歴情報 1 2 3 をマスターサーバ 2 0 0 にそれぞれ送信する。

【 0 0 6 2 】

各サーバ 1 0 0 は、マスターサーバ 2 0 0 に自サーバ内アクセス履歴情報 1 2 3 を送信し、その後、リード回数 1 2 3 3 及びライト回数 1 2 3 4 の値を変更する。変更方法とし

10

20

30

40

50

ては、例えば回数を0にする、半分にする、などの値を減ずる方法を取る。

【0063】

マスターサーバ200は、ステップS403において、各サーバ100より自サーバ内アクセス履歴情報123をそれぞれ受信し、これらの自サーバ内アクセス履歴情報123を元に、クラスタ内アクセス履歴情報223を更新する。

【0064】

これにより、マスターサーバ200ではアクセス履歴取得単位2231毎に、全サーバ100のアクセス回数2233、2234の合計2235を算出する。この時、サーバ100毎に重み付けを行ったり、リードとライトで異なった重み付けを行う、等の重み付けを行うようにしても良い。あるサーバ100に対するアクセス履歴取得単位2231の重み付けを無限大に設定することにより、データをそのサーバ100に固定することも可能である。

10

【0065】

次に、マスターサーバ200は、ステップS404において、クラスタ内アクセス履歴情報223のアクセス履歴取得単位2231毎に、全サーバ100からのアクセスの合計2235が高い順に並び替える。この並び替えは、サーバ合計2235のリード回数2235Rとライト回数2235Wの和や、リード回数2235Rとライト回数2235Wの一方など、予め設定した基準で行えば良い。

【0066】

そして、マスターサーバ200は、アクセス履歴取得単位2231のアクセス数の合計2235が所定の閾値以上のものに対し、アクセス回数が高い順にアクセス履歴取得単位2231を並べ替える。そして、サーバ100の不揮発性メモリ140の容量を超えない量のデータを順次配置する。なお、上記閾値は、サーバ合計2235のリード回数2235Rとライト回数2235Wの和に対して設定する閾値や、リード回数2235Rとライト回数2235Wのそれぞれに対する閾値など、所定の基準で設定することができる。

20

【0067】

そしてステップS405において、マスターサーバ200は、サーバ100毎のアクセス履歴取得単位2231へのアクセス回数2233、2234に応じて、どのデータをどのサーバ100-1~100-nに配置するかを決定する。

【0068】

この配置の決定は、ステップS404でマスターサーバ200が不揮発性メモリ140上に配置すると決定したデータそれぞれに対し、そのデータへのアクセス回数が多いサーバ100-1~100-nの順に、不揮発性メモリ140への配置を試みる。このとき、マスターサーバ200は、サーバ100の不揮発性メモリ140に容量が余っていればそのサーバ100に、ステップS404で配置されたデータを格納し、不揮発性メモリ140の容量が足りなければ、そのサーバ100の不揮発性メモリ140の余っている容量分だけデータを配置し、入らなかったデータは次にアクセス回数が多いサーバ100に配置するものとする。なお、アクセス回数と同じサーバ100が複数存在する場合は、どのサーバ100に配置するかを、例えば、余っている不揮発性メモリ140容量のより大きなサーバ100に決定する、ランダムに決定する、等の周知または公知の方法を取ることができる。

30

40

【0069】

図11は、図9のステップS301に示す処理で、マスターサーバ200が新たなデータの配置に従って、各サーバ100に不揮発性メモリ140へのデータ登録または削除を指示する処理の詳細なフローチャートを示す。

【0070】

まず、マスターサーバ200は、ステップS501において、不揮発性メモリ140からの削除対象データを抽出し、それぞれのデータ削除の処理を行う。この削除データの抽出処理は、例えば、アクセス回数の低下によって、不揮発性メモリ140から削除し、共有ストレージシステム500のみでデータの読み書きを行う配置になったデータを、マ

50

ターサーバ200が不揮発性メモリ140からの削除対象として抽出する。そして、マスターサーバ200は、現在不揮発性メモリ140に格納している移動対象のデータを削除する指令をサーバ100に通知する。そして、この通知を受信したサーバ100は、通知されたデータを不揮発性メモリ140から削除する。なお、削除処理の詳細については、図13において後述する。

【0071】

次に、ステップS502において、マスターサーバ200はあるサーバ100の不揮発性メモリ140から別のサーバ100-nの不揮発性メモリ140へ移すデータを抽出し、それぞれのデータに対し、まず上記ステップS501同様のデータの削除処理を行った後、データの登録処理を行う。登録処理は、マスターサーバ200が、移動先のサーバ100-nに、登録するデータの共有ストレージシステム500上のセクタ番号を通知し、当該サーバ100にデータの登録を指令する。通知を受信したサーバ100は、指定されたセクタ番号のデータを、ストレージシステム500から読み出して不揮発性メモリ140に格納する。なお、登録処理の詳細については、図12において後述する。

10

【0072】

最後に、ステップS503において、マスターサーバ200は不揮発性メモリ140へ新たに登録対象となったデータに対し、それぞれのデータの追加の処理を行う。新たに登録対象となったデータは、現在不揮発性メモリ140には格納されておらず、共有ストレージシステム500のみに格納されているが、アクセス回数の増大などで、不揮発性メモリ140へ登録することが決定されたデータを示す。

20

【0073】

ここで、上記ステップS502の処理において、全てのデータの削除処理が完了する前に登録処理が行われる可能性がある為、あるサーバ100の不揮発性メモリ140の容量が不足する可能性がある。この場合は、別データの移動処理を先に処理すれば、容量の足りないサーバ100の不揮発性メモリ140への削除処理がいずれかのタイミングで行われるので、問題が解決する。

【0074】

また、上記ステップS502の処理で削除処理を登録処理より先に行うのは、サーバ100間でのコヒーレンシを保つ為である。すなわち、登録処理を先に行った場合、システム上にあるデータのコピーが複数存在することとなり、この状態でいずれかのサーバ100がライトを行った場合、複数存在する全てのデータへ同時に書き込まなければコヒーレンシが失われる可能性がある。しかしながら、本計算機システムのライトは図8に示したように行われるため、このような仕組みを有していない。そこで、データの削除を行ってからデータの登録を行うことによって、不揮発性メモリ140上に存在する共有ストレージシステム500のデータのコピーの数を1つに保ち、コヒーレンシを保つようにする。

30

【0075】

図12は、図11のステップS502、S503における各サーバ100の不揮発性メモリ140へのデータ登録処理の詳細なフローチャートを示す。

【0076】

各サーバ100の不揮発性メモリ140へのデータ登録処理において、マスターサーバ200はまず、ステップS601において、データの登録対象のサーバ100へデータの登録要請を行い、登録するデータの共有ストレージセクタ番号を通知する。

40

【0077】

ここで、データ登録先のサーバ100の不揮発性メモリ140の容量は、マスターサーバ200側のクラスタ内不揮発性メモリ利用情報222の不揮発性メモリ合計容量222と、利用容量2223で保持しているため、データ登録先のサーバ100の不揮発性メモリ140の容量不足でデータ登録が行えないという問題は発生しない。

【0078】

次に、ステップS602において、データ登録元のサーバ100は、自サーバ内不揮発性メモリ利用情報122を参照して、自サーバ内の不揮発性メモリ140のデータ登録を

50

行うデバイス番号 1 2 2 1 及び不揮発性メモリセクタ番号 1 2 2 4 を決定する。

【 0 0 7 9 】

そして、ステップ S 6 0 3 において、データ登録元のサーバ 1 0 0 はマスターサーバ 2 0 0 へデータ登録位置を通知する。このデータ登録位置は、不揮発性メモリ格納位置情報 1 2 4 のデータ格納位置 1 2 4 2 と同様に、サーバ 1 0 0 の識別子 (サーバ番号 2 2 2 1) と、デバイス番号、不揮発性メモリセクタ番号を含む。

【 0 0 8 0 】

次に、ステップ S 6 0 4 において、マスターサーバ 2 0 0 は全サーバ 1 0 0 に対し、不揮発性メモリ格納位置情報 1 2 4 の共有ストレージセクタ番号 1 2 4 1 がデータ登録対象であるエントリに対し、RDMA 利用禁止モードへの変更を依頼する。なお、RDMA の利用禁止モードは、各サーバ 1 0 0 の不揮発性メモリ格納位置情報 1 2 4 の RDMA 利用可否 1 2 4 3 を空白にすればよい。

10

【 0 0 8 1 】

各サーバ 1 0 0 はステップ S 6 0 5 において、RDMA が利用出来ないモードで不揮発性メモリ格納位置情報 1 2 4 を更新し、マスターサーバ 2 0 0 へレスポンスを返す。マスターサーバ 2 0 0 は、ステップ S 6 0 6 において、全てのサーバ 1 0 0 からのレスポンスを受け取り、不揮発性メモリ格納位置情報 1 2 4 の更新完了をデータ登録元のサーバ 1 0 0 に通知する。

【 0 0 8 2 】

そして、ステップ S 6 0 7 においてデータ登録元のサーバ 1 0 0 は、共有ストレージシステム 5 0 0 より登録対象のデータを読み込み、ステップ S 6 0 8 において、共有ストレージシステム 5 0 0 より読み出したデータを不揮発性メモリ 1 4 0 へ書き込む。

20

【 0 0 8 3 】

そして、データ登録元のサーバ 1 0 0 は、ステップ S 6 0 9 においてマスターサーバ 2 0 0 へ登録完了を通知する。マスターサーバ 2 0 0 はステップ S 6 1 0 において、データ登録先の不揮発性メモリ 1 4 0 で RDMA が使用可能であった場合は、全てのサーバ 1 0 0 の不揮発性メモリ格納位置情報 1 2 4 の該当エントリに対して、RDMA の利用が可能なモードへの変更を通知する。そして最後に、ステップ S 6 1 1 において、全てのサーバ 1 0 0 は不揮発性メモリ格納位置情報 1 2 4 を更新する。

【 0 0 8 4 】

30

ここで、ステップ S 6 0 3 ~ S 6 0 6 までの処理は、登録するデータの coherence 制御の為に、すなわち、データ登録元のサーバ 1 0 0 が、ある時点で共有ストレージシステム 5 0 0 のデータを不揮発性メモリ 1 4 0 に取り込んだ場合、データ取り込み中に別のサーバ 1 0 0 - n がそのデータを更新する可能性がある。この場合、不揮発性メモリ 1 4 0 上のデータと共有ストレージシステム 5 0 0 上のデータが異なる状態となり、coherence が失われる。そこで、最初に全てのサーバ 1 0 0 の不揮発性メモリ格納位置情報 1 2 4 をデータの更新を行うサーバ 1 0 0 に RDMA の利用を禁止することにより、データを登録するサーバ 1 0 0 にデータ取り込み開始から終了までの全てのデータ更新を把握させる。

【 0 0 8 5 】

40

データを取り込むサーバ 1 0 0 は、データ取り込み終了前に各サーバ 1 0 0 から行われたリード若しくはライトに対し、リードは共有ストレージシステム 5 0 0 から読み出し、ライトはデータ登録元のサーバ 1 0 0 のバッファ内に記録しておく。そして、不揮発性メモリ 1 4 0 へバッファのデータを登録する際、すなわちステップ S 6 0 8 の処理において、共有ストレージシステム 5 0 0 から読み出したデータの上書きを行う。これによって、共有ストレージシステム 5 0 0 とデータ登録元のサーバ 1 0 0 の不揮発性メモリ 1 4 0 上のデータの coherence が保たれる。RDMA を利用することにより I/O 性能を向上可能なので、最後にステップ S 6 1 0、S 6 1 1 で RDMA を利用可能なモードへ変更を行っている。

【 0 0 8 6 】

50

図13は、図11のステップS501、S502におけるマスターサーバ200と各サーバ100の不揮発性メモリ140のデータ削除処理の詳細なフローチャートを示す。

【0087】

各サーバ100の不揮発性メモリ140へのデータ削除処理において、マスターサーバ200はまず、ステップS701において、削除対象のデータを保持するサーバ100以外の全てのサーバ100へ、不揮発性メモリ格納位置情報124の該当エントリの削除を依頼する。

【0088】

各サーバ100はステップS702において、不揮発性メモリ格納位置情報124の該当エントリを削除後、エントリを削除する前の全てのI/Oのレスポンスを受信してからマスターサーバ200へ応答を返す。

10

【0089】

ステップS703でマスターサーバ200は、全てのサーバ100からのレスポンスの到着を待って、その後、ステップS704において、マスターサーバ200は削除対象データを保持するサーバ100へデータの削除を依頼する。

【0090】

そして、ステップS705において、データを削除するサーバ100は不揮発性メモリ格納位置情報124の該当エントリを削除し、最後にステップS706において、マスターサーバ200に削除完了を通知する。

【0091】

ここで、ステップS702において、マスターサーバ200へすぐ応答を返さず、削除前の全てのI/O処理のレスポンスを受信後にマスターサーバ200へレスポンスを返すのは、該当エントリを削除する前のI/Oが遅延した場合、宛先に新しいデータが書き込まれてからリード・ライトが行われるのを防ぐためである。

20

【0092】

以上、図7～図13の処理によって、サーバ100に搭載された不揮発性メモリ140と共有ストレージシステム500間のデータの階層化制御と、サーバ100間のデータ配置の最適化制御を動的に行うことが達成される。

【0093】

図14は、マスターサーバ200及びサーバ100で行われる初期化のフローチャートを示す。初期化は、図2～図6に示す表の初期化として表される。まず、ステップS801において、マスターサーバ200にどのサーバ100の不揮発性メモリ140のどれだけの容量をサーバ100間で共有するかを設定する。

30

【0094】

これは、例えば、マスターサーバ200の図示しない入力装置に対して人手で容量を入力する方法や、マスターサーバ200が各サーバ100より自動的に情報を収集し、容量の設定を行う方法が考えられる。

【0095】

次に、ステップS802において、マスターサーバ200において必要に応じて人手でデータの初期配置を指定する。例えば、共有ストレージシステム500のデータの内、頻繁にアクセスされるデータが予め分かっている場合は、予め各サーバ100の不揮発性メモリ140に配置するデータを管理者等が設定しておく。これにより、システム起動時より不揮発性メモリ140の搭載による性能向上を図る事ができる。

40

【0096】

そして、マスターサーバ200はステップS803において、クラスタ内不揮発性メモリ利用情報222のうち、各サーバ100の不揮発性メモリ140の容量をステップS801で収集した情報から不揮発性メモリ合計容量2222を設定し、ステップS802の設定状況に応じて利用容量2223と、不揮発性メモリ保持共有ストレージセクタ番号2224を設定する。次に、マスターサーバ200はステップS804において、各サーバ100にクラスタ内不揮発性メモリ利用情報222及びS802のデータ初期配置を配布

50

する。そして各サーバ100は、クラスタ内不揮発性メモリ利用情報222に従って、自サーバ内不揮発性メモリ利用情報122の容量1222を設定する。また各サーバ100は、自サーバ内不揮発性メモリ利用情報122の利用容量1223、不揮発性メモリセクタ番号1224及び対応共有ストレージセクタ番号1225を、S802のデータ初期配置に応じて設定を行う。

【0097】

次に、ステップS805において、マスターサーバ200のクラスタ内アクセス履歴情報223にアクセス履歴取得単位2231を設定する。これは、200Mbyte毎等の一定の大きさの単位で切る方法や、各サーバ100で稼働するアプリケーションが利用するデータの種別毎に切る方法が考えられる。後者の例としては、データベースを扱う場合に、インデックスとデータを別の単位として設定する事が考えられる。そして、ここで設定したアクセス履歴取得単位2231に応じてクラスタ内アクセス履歴情報223のアクセス履歴取得単位2231及び共有ストレージセクタ番号2232を設定する。マスターサーバ200は、クラスタ内アクセス履歴情報223のアクセス回数2233、2234は全て0に設定する。

10

【0098】

マスターサーバ200はステップS806において、各サーバ100にクラスタ内アクセス履歴情報223を配布する。各サーバ100は配布された情報に合わせて自サーバ内アクセス履歴情報123の共有ストレージセクタ番号1232を設定し、アクセス回数1233、1234は全て0に設定する。

20

【0099】

最後にステップS807において、マスターサーバ200を含めた全てのサーバ100の不揮発性メモリ格納位置情報124、224に対して、全ての共有ストレージセクタ番号のデータ格納位置1242を共有ストレージシステム500に設定する。

【0100】

以上の処理により、全てのテーブルの初期化が完了する。

【0101】

このような環境では、あるサーバ100の電源が遮断された場合、他のサーバ100-nがそのサーバ100の不揮発性メモリ140へアクセスを行ってしまうと、レスポンスが何時まで経っても帰ってこない状況が発生する、という問題がある。その為、サーバ100の電源遮断を行うときには、まずマスターサーバ200が電源遮断対象のサーバ100の不揮発性メモリ140を利用しないように再設定を行い、その後サーバ100の電源遮断を行う必要がある。この時の処理の一例を示すフローチャートを図15に示す。

30

【0102】

図15は、電源遮断時にマスターサーバ及びサーバで行われる処理の一例を説明するフローチャートである。

【0103】

まず、ステップS901において、電源を遮断するサーバ100はマスターサーバ200に電源遮断の通知を行う。

【0104】

マスターサーバ200は電源遮断の通知をサーバ100から受け取った後、ステップS902において、電源遮断対象のサーバ100の不揮発性メモリ140上にあるデータの削除処理を実施する。これにより、全サーバ100は電源遮断対象のサーバ100の不揮発性メモリ140へアクセスしなくなる。

40

【0105】

マスターサーバ200は、ステップS903において、電源遮断対象のサーバ100に削除処理の完了を通知する。これにより、電源遮断対象のサーバ100は通常の電源遮断処理に戻る。

【0106】

次にマスターサーバ200は、ステップS904において、クラスタ内アクセス履歴情

50

報 2 2 3 より電源遮断対象のサーバ 1 0 0 以外のデータ配置に関し、不揮発性メモリ 1 4 0 から削除するデータと、新たに登録するデータを決定する。この処理は、電源遮断対象のサーバ 1 0 0 上のデータを削除したため、一時的に配置が最適化されていない状態になっている為、実施する。最後に、マスターサーバ 2 0 0 はステップ S 9 0 5 において、データの削除処理と登録処理を行う。

【 0 1 0 7 】

なお、サーバ 1 0 0 のクラスタ環境において、どのサーバ 1 0 0 にどのデータを配置するかを静的に人手で決定することも可能である。しかし、クラスタ環境は、動作するアプリケーションが昼と夜で全く異なる事や、数ヶ月～数年に一度はシステムのリプレイス若しくは増強が行われ、サーバ 1 0 0 の台数やサーバ 1 0 0 毎の不揮発性メモリ 1 4 0 の搭載容量が変化することを考えると、静的に人手でデータ配置を考えるのはほぼ不可能であると考えられる。本計算機システムでは動的にデータ配置をマスターサーバ 2 0 0 が決定するため、上で述べた状況にも対応可能である。

10

【 0 1 0 8 】

上記実施例においては、共有ストレージシステム 5 0 0 上のデータの位置情報として、セクタ番号を用いる例を示したが、ブロック番号や論理ブロックアドレスなどを用いるようにしても良い。

【 0 1 0 9 】

なお、本発明において説明したサーバ等の構成、処理部及び処理手段等は、それらの一部又は全部を、専用のハードウェアによって実現してもよい。

20

【 0 1 1 0 】

また、本実施例で例示した種々のソフトウェアは、電磁的、電子的及び光学式等の種々の記録媒体（例えば、非一時的な記憶媒体）に格納可能であり、インターネット等の通信網を通じて、コンピュータにダウンロード可能である。

【 0 1 1 1 】

また、本発明は上記した実施例に限定されるものではなく、様々な変形例が含まれる。例えば、上記した実施例は本発明をわかりやすく説明するために詳細に説明したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。

【 符号の説明 】

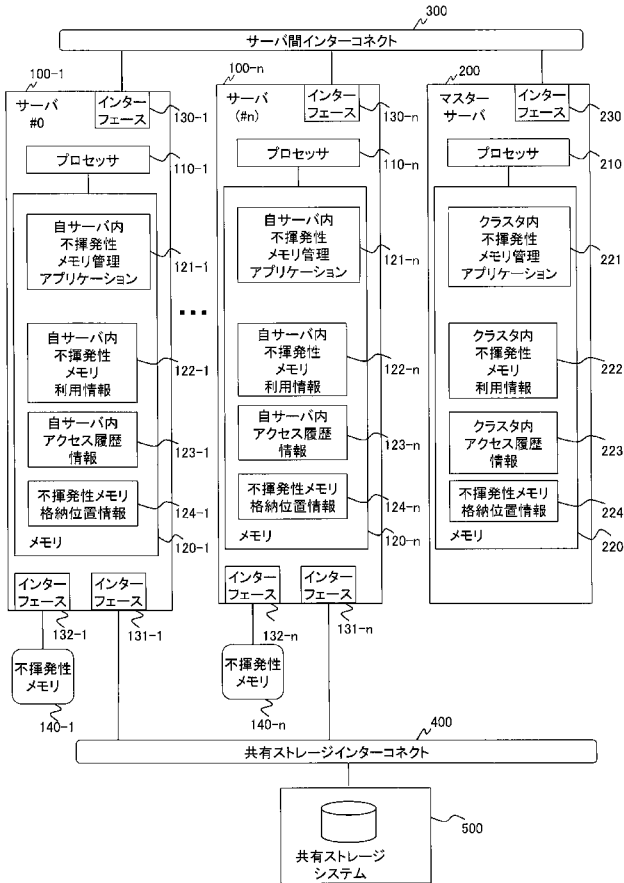
【 0 1 1 2 】

- 1 0 0 - 1 ~ 1 0 0 - n サーバ
- 1 1 0 - 1 ~ 1 1 0 - n、2 1 0 プロセッサ
- 1 2 0 - 1 ~ 1 2 0 - n メモリ
- 1 2 2 - 1 ~ 1 1 2 - n 自サーバ不揮発性メモリ利用情報
- 2 2 2 クラスタ内不揮発性メモリ利用情報
- 1 2 3 - 1 ~ 1 2 3 - n 自サーバ内アクセス履歴情報
- 2 2 3 クラスタ内アクセス履歴情報
- 1 2 4 - 1 ~ 1 2 4 - n、2 2 4 不揮発性メモリ格納位置情報
- 1 4 0 不揮発性メモリ
- 2 0 0 マスターサーバ
- 3 0 0 サーバ間インターコネク
- 4 0 0 共有ストレージインターコネク
- 5 0 0 共有ストレージシステム

30

40

【 図 1 】



【 図 2 】

122 自サーバ内不揮発性メモリ利用情報

1221 不揮発性メモリデバイス番号	1222 容量	1223 利用容量	1224 不揮発性メモリセクタ番号	1225 対応共有ストレージセクタ番号
0	400GB	200GB	0 - 419,430,400	0 - 419,430,400
			419,430,400 - 838,860,800	利用なし
1	200GB	200GB	0 - 419,430,400	838,860,800 - 1,258,291,200

【 図 3 】

222 クラスタ内不揮発性メモリ利用情報

2221 サーバ番号	2222 不揮発性メモリ合計容量	2223 利用容量	2224 不揮発性メモリ内保持共有ストレージセクタ番号
サーバ#0	600GB	400GB	0 - 419,430,400 838,860,800 - 1,258,291,200
サーバ#1	500GB	200GB	419,430,400 - 838,860,800
...			

【 図 4 】

124, 224 不揮発性メモリ格納位置情報

1241 共有ストレージセクタ番号	1242 データ格納位置	1243 RDMA利用可否
0 - 419,430,400	サーバ#0 Vol#0 0 - 419,430,400	○
419,430,400 - 838,860,800	サーバ#1 Vol#0 838,860,800 - 1,258,291,200	○
838,860,800 - 1,258,291,200	サーバ#0 Vol#1 0 - 419,430,400	○
その他	共有ストレージ	

【 図 6 】

223 クラスタ内アクセス履歴情報

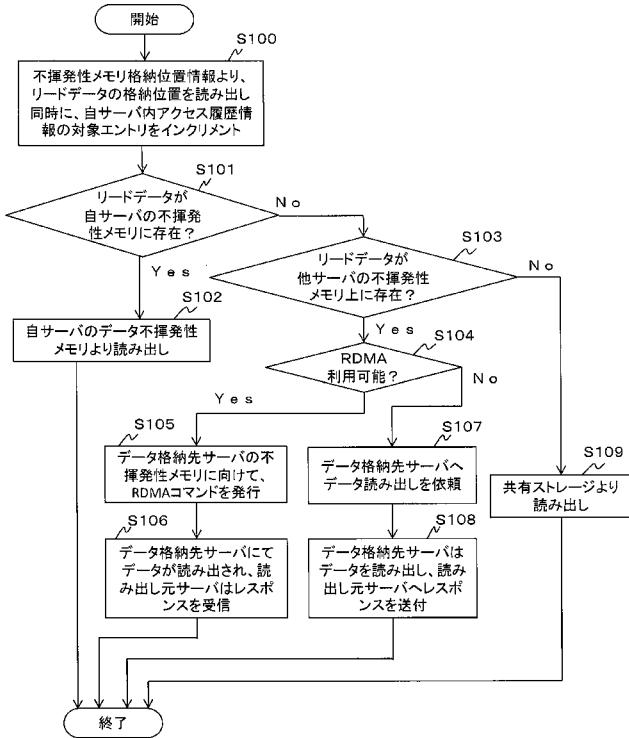
2231 アクセス履歴取得単位	2232 共有ストレージセクタ番号	2233 サーバ#0		2234 サーバ#1		2235 サーバ合計	
		リード回数	ライト回数	リード回数	ライト回数	リード回数	ライト回数
1	0 - 20,480	153	250	50	1520	203	1770
2	20,480 - 40,960	2,103	125	1,653	0	3,756	125
3	40,960 - 61,440	12,453	0	53,320	0	65,773	0
...	...						

【 図 5 】

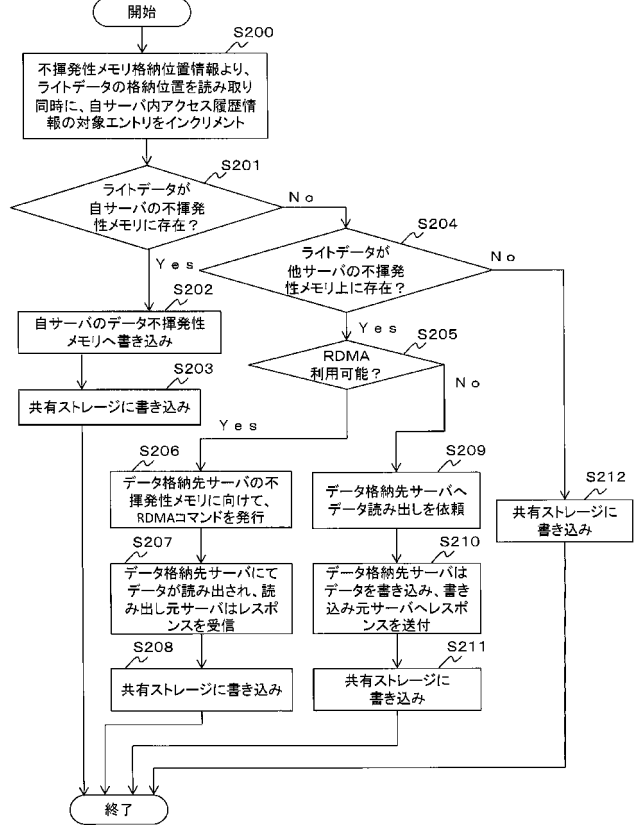
123 自サーバ内アクセス履歴情報

1231 アクセス履歴取得単位	1232 共有ストレージセクタ番号	1233 リード回数	1234 ライト回数
1	0 - 20,480	153	250
2	20,480 - 40,960	0	0
3	40,960 - 61,440	12453	0
...	...		

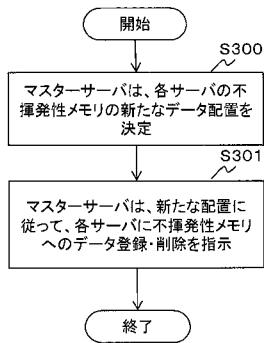
【 図 7 】



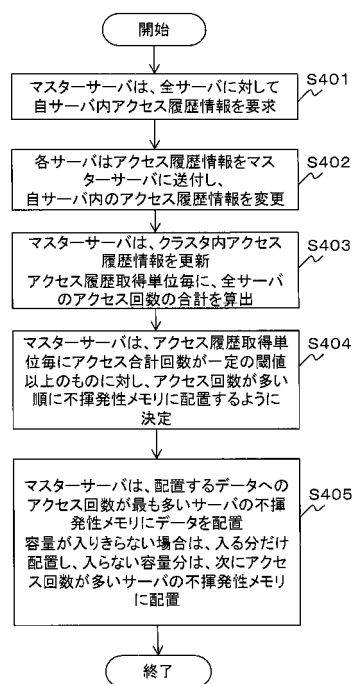
【 図 8 】



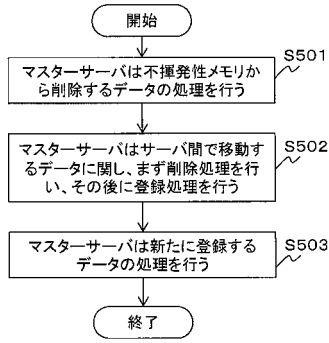
【 図 9 】



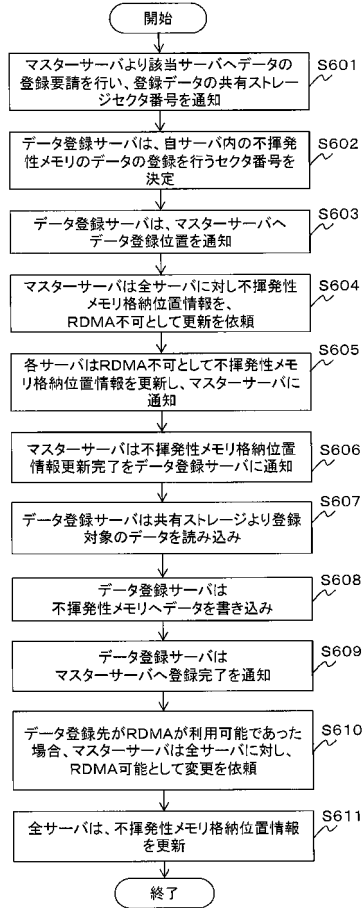
【 図 10 】



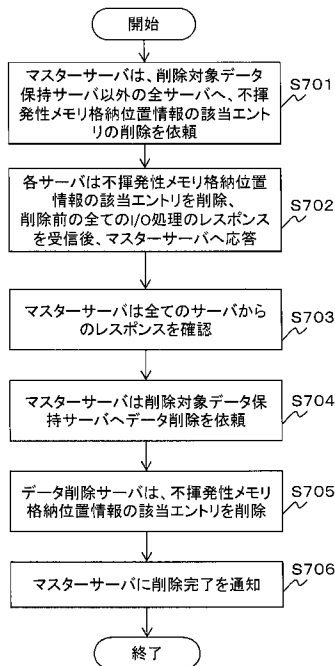
【図 1 1】



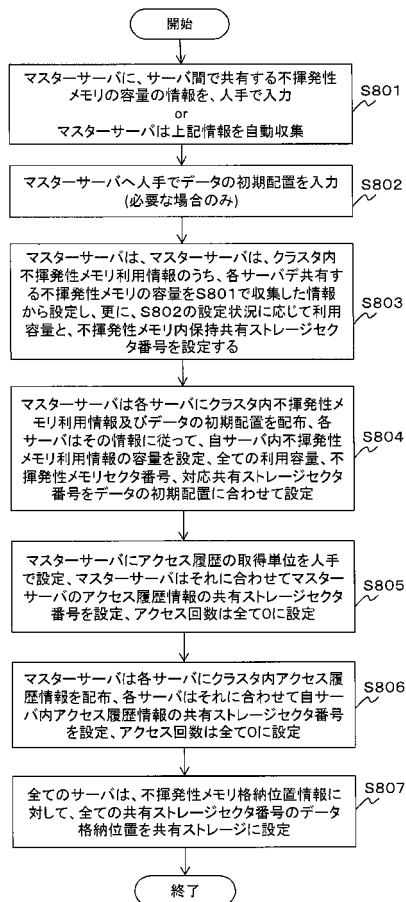
【図 1 2】



【図 1 3】



【図 1 4】



【 図 1 5 】

