

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3704573号

(P3704573)

(45) 発行日 平成17年10月12日(2005.10.12)

(24) 登録日 平成17年8月5日(2005.8.5)

(51) Int. Cl.⁷

G06F 12/00

F I

G06F 12/00 535Z

G06F 12/00 514E

G06F 12/00 545A

請求項の数 3 (全 15 頁)

(21) 出願番号	特願2001-71680 (P2001-71680)	(73) 特許権者	301063496
(22) 出願日	平成13年3月14日 (2001.3.14)		東芝ソリューション株式会社
(65) 公開番号	特開2002-268933 (P2002-268933A)		東京都港区芝浦一丁目1番1号
(43) 公開日	平成14年9月20日 (2002.9.20)	(74) 代理人	100058479
審査請求日	平成13年3月14日 (2001.3.14)		弁理士 鈴江 武彦
		(74) 代理人	100091351
			弁理士 河野 哲
(出願人による申告) 国などの委託研究の成果に係る特許出願 (平成12年度通信・放送機構「スーパーインターネットプラットフォーム技術の研究開発」委託研究、産業活力再生特別措置法第30条の適用を受けるもの)		(74) 代理人	100088683
			弁理士 中村 誠
		(74) 代理人	100108855
			弁理士 蔵田 昌俊
		(74) 代理人	100075672
			弁理士 峰 隆司
		(74) 代理人	100109830
			弁理士 福原 淑弘

最終頁に続く

(54) 【発明の名称】 クラスタシステム

(57) 【特許請求の範囲】

【請求項1】

複数の電子計算機と、共有ディスク装置と、前記複数の電子計算機上で動作するプロセスから前記共有ディスク装置に記録されたファイルに対して共有アクセスするためにロック機能によりデータの一貫性を保持するための排他制御をするクラスタ共有ファイルシステムとを持ったクラスタシステムであって、

前記電子計算機は、

前記クラスタ共有ファイルシステムが管理するファイルを前記プロセスのアドレス空間内にマッピングして共有メモリを設定するためのクラスタ共有メモリ設定手段と、

前記クラスタ共有ファイルシステムのロック機能を前記共有メモリにアロケーションして前記共有メモリ上でデータの一貫性を保持するための排他制御を可能とするクラスタ共有メモリ用ロックアロケーション手段と、

前記設定された共有メモリ内の全てのページに対するアクセスを禁止にするアクセス禁止設定手段と、

前記アクセス禁止が設定されたページをアクセスしたときにREADページフォールトが発生した場合、そのページに前記マッピングされたデータを前記共有ディスク装置に記録されたファイルから読み出して書き込むデータ書き込み手段と、

前記データが書き込まれたページを読み出し可能に設定する設定手段とを具備することを特徴とするクラスタシステム。

【請求項2】

10

20

複数の電子計算機と、共有ディスク装置と、前記複数の電子計算機上で動作するプロセスから前記共有ディスク装置に記録されたファイルに対して共有アクセスするためにロック機能によりデータの一貫性を保持するための排他制御をするクラスタ共有ファイルシステムとを持ったクラスタシステムであって、

前記電子計算機は、

前記クラスタ共有ファイルシステムが管理するファイルを前記プロセスのアドレス空間内にマッピングして共有メモリを設定するためのクラスタ共有メモリ設定手段と、

前記クラスタ共有ファイルシステムのロック機能を前記共有メモリにアロケーションして前記共有メモリ上でデータの一貫性を保持するための排他制御を可能とするクラスタ共有メモリ用ロックアロケーション手段と、

前記設定された共有メモリ内の全てのページに対するアクセスを禁止にするアクセス禁止設定手段と、

前記アクセス禁止が設定されたページをアクセスしたときにWRITEページフォールトが発生した場合、そのページに前記マッピングされたデータを前記共有ディスク装置に記録されたファイルから読み出して書き込むデータ書き込み手段と、

前記データが書き込まれたページを読み出し/書き込み可能に設定する設定手段とを具備することを特徴とするクラスタシステム。

【請求項3】

前記クラスタ共有メモリをロック操作した場合、前記クラスタ共有ファイルシステムが管理するファイルを前記プロセスのアドレス空間にマッピングしたクラスタ共有メモリ領域のページをアクセス不許可に設定してロックを取得する手段と、

前記取得したロックをアンロック操作した場合、前記クラスタ共有メモリ領域の更新されたページのデータを前記クラスタ共有ファイルシステムが管理するファイルに書き戻す手段を具備したことを特徴とする請求項2記載のクラスタシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、記憶装置と計算機を専用の高速ネットワークで接続したSAN(Storage Area Network(ストレージエリア・ネットワーク))環境における並列処理プログラミングが容易なクラスタシステムに関するものである。

【0002】

【従来の技術】

近年、データのストレージ装置としての磁気ディスク装置と電子計算機(以下、サーバーコンピュータまたは単にサーバーと呼ぶ)とをFiber Channel(ファイバ・チャンネル)等の専用の高速ネットワークで接続したSANが注目を浴びている。

【0003】

このSANには以下のような利点がある。(1)複数のサーバーがストレージ装置を共有することができる。(2)ストレージ装置のアクセス負荷をLAN(Local Area Network)から分離させることができる。(3)ファイバ・チャンネル等によりストレージ装置へのアクセスを高速にすることができる。このうち(1)はSANの一般的な利点であるが、ここにクラスタ共有ファイルシステムの技術を用いると、単に複数のサーバーがストレージ装置を共有できるだけでなく、ファイルを共有アクセスすることが可能になる。

【0004】

米国Sistina Software Incが開発しているGFS(Global File System)は、このようなクラスタ共有ファイルシステムの一例である。クラスタ共有ファイルシステムでは、複数のサーバーコンピュータによるストレージ装置に記憶されたファイルへの共有アクセスを可能にしている。

【0005】

複数のサーバーコンピュータによるファイルへの共有アクセスというと、一般的にはNF

10

20

30

40

50

S (Network File System) が連想されるが、NFSでは複数のサーバーコンピュータで同一ファイルを更新した場合のデータの一貫性を保証していない。これに対しクラスタ共有ファイルシステムではデータの一貫性を保証している。

【 0 0 0 6 】

クラスタ共有ファイルシステムは、複数のサーバーコンピュータによるファイルの共有アクセス (READ / WRITE) 機能の他に、データの一貫性を保証するために複数のサーバーコンピュータに跨るロック機構であるクラスタ共有ファイルシステム用ロック機構を持っている。このクラスタ共有ファイルシステム用ロック機構により、複数のサーバーコンピュータから構成されるクラスタシステム (Cluster System) 上で並列アプリケーションプログラムが実行可能となる。

10

【 0 0 0 7 】

例えば、クラスタ共有ファイルシステムとクラスタ共有ファイルシステム用ロック機構とを実装した複数のサーバーコンピュータがSANを介して1つの磁気ディスク装置に接続されているクラスタシステムについて考察する。このクラスタシステムは、主メモリを共有していない疎結合のクラスタシステムである。各サーバーコンピュータ上で実行されているプロセスは、クラスタ共有ファイルシステムを用いることにより、磁気ディスク装置に記憶されているファイルに共有アクセスすることができる。しかも、プロセスは、クラスタ共有ファイルシステム用ロック機構を用いて排他制御処理をすることにより、磁気ディスク装置に記憶されているファイルのデータの一貫性を保証することができる。

20

【 0 0 0 8 】

これに対して、ファイルシステムと共有メモリシステムとロック機構とが実装され、複数のプロセッサを持ち、例えば1つの磁気ディスク装置が接続されているSMP型 (Symmetrical Multiprocessor) 並列計算機について考察する。このSMP型並列計算機上で実行されている複数のプロセスは、それぞれファイルシステムを介して磁気ディスク装置に記憶されたファイルに共有アクセスしたり、共有メモリシステムを介して共有メモリ (主メモリ) にアクセスすることができる。また、複数のプロセスは、ロック機構を介して磁気ディスク装置に記憶されたファイルや共有メモリに記憶されたデータに対する共有アクセスにおける排他制御処理をすることにより、データの一貫性を保持することができる。

30

【 0 0 0 9 】

このように従来の疎結合のクラスタシステムとSMP型並列計算機とを比較すると、両者ともに磁気ディスク装置に記憶されたファイルに対する共有アクセスおよびこれらに対するデータの一貫性の保持は可能である。しかし、クラスタシステムでは、複数のサーバーコンピュータから共有メモリへの共有アクセスをすることができない。

【 0 0 1 0 】

【 発明が解決しようとする課題 】

クラスタ共有ファイルシステムを持つクラスタシステムでは、ファイルへの共有アクセスはできるが、メモリへの共有アクセスはできない。このため、クラスタ共有ファイルシステムを持つクラスタシステムでは、SMP型並列計算機と比べると、その実行するアプリケーションプログラムを並列プログラムにて記述するのが困難であるという課題がある。具体的には、複数のサーバーコンピュータ上で実行されているプロセスでファイル (又はデータ) を共有して互いにデータの一貫性を確保しながらデータの同期をとってデータ処理するプログラミングを記述する場合には、プロセスが共有して処理するファイル (又はデータ) をメモリ上に配置して処理することが必要であるが、これを行うことができない。もし、プロセスがファイル (又はデータ) を同期をとりながらデータの一貫性を確保して共有処理することが必要な場合には、クラスタ共有ファイルシステムを利用して、その共有対象のファイル (又はデータ) を磁気ディスク装置上に配置して処理することが必要である。この場合には、プロセスが、ファイル (又はデータ) に対する処理を入出力装置を対象とするI/O処理のコマンドを使用してプログラミングされている必要がある。こ

40

50

のようにファイル装置（磁気ディスク装置）に記憶されているファイル（又はデータ）のデータ処理する場合は、メモリ上に配置されたファイル（又はデータ）を処理する場合に比べて、プログラムの記述が複雑になる。これは、メモリ上に配置されたデータに対する処理であれば、load命令やstore命令を使用して簡単なプログラム記述で処理が記述できる。

【0011】

しかし、これに対して磁気ディスク装置上に配置されたファイルやデータを処理する場合には、I/O処理をするための複雑なコマンドを使用したプログラムを記述する必要があり、並列プログラムにて記述するのが困難である。また、メモリ上に配置されたデータに対するデータ処理に比べて、磁気ディスク装置などのファイル装置（I/O装置）上に配置されたデータを処理する場合には、その処理時間が多くかかり、処理速度が落ちるとい

10

【0012】

以上説明したように、クラスタ共有ファイルシステムを実装した複数のサーバーコンピュータが疎結合されたクラスタシステムでは、磁気ディスク装置などのファイル装置上に記憶されたファイルへの共有アクセスは可能であるが、主メモリへの共有アクセスは可能でないため、SMP型並列計算機と比べると、並列プログラムにて記述するのが困難であるという課題がある。

【0013】

本発明は、クラスタ共有ファイルシステムを実装した複数のサーバーコンピュータが疎結合されたクラスタシステムにおいて、並列プログラムの記述が容易なクラスタシステムを提供することを目的とする。

20

【0014】**【課題を解決するための手段】**

本発明は、複数の電子計算機と、共有ディスク装置と、前記複数の電子計算機上で動作するプロセスから前記共有ディスク装置に記録されたファイルに対して共有アクセスするためにロック機能によりデータの一貫性を保持するための排他制御をするクラスタ共有ファイルシステムとを持ったクラスタシステムであって、前記クラスタ共有ファイルシステムが管理するファイルを前記プロセスのアドレス空間内にマッピングして共有メモリを設定するためのクラスタ共有メモリ設定手段と、前記クラスタ共有ファイルシステムのロック機能を前記共有メモリにアロケーションして前記共有メモリ上でデータの一貫性を保持するための排他制御を可能とするクラスタ共有メモリ用ロックアロケーション手段とを備えたことを特徴とする。

30

【0015】

本発明によれば、クラスタ共有ファイルシステムを実装した複数のサーバーコンピュータが疎結合されたクラスタシステムにおいて、並列プログラムの記述を容易にすることができる。

【0016】**【発明の実施の形態】**

本発明のポイントは、クラスタシステムにおいて、アプリケーションプログラムを実行することで生成されるプロセスが共有ディスク装置に記録されたファイルを、クラスタ共有ファイルシステムを用いてプロセスが主メモリ上に配置されているアドレス空間内に仮想的に設けたクラスタ共有メモリ（分散共有メモリ）領域上にマッピングすることで、ファイルへのアクセスを主メモリへのアクセスとして処理することである。

40

【0017】

以下、図面を用いて、本発明の一実施形態を詳細に説明する。図1は、本発明のクラスタシステムを示す図である。図1において、クラスタシステム10は、ファイバ・チャンネルで構成されているストレージエリア・ネットワーク（SAN）15にサーバーコンピュータ11、サーバーコンピュータ12、サーバーコンピュータ13、サーバーコンピュータ14と磁気ディスク装置で構成される共有ディスク装置16が接続されて構成されている

50

。

【0018】

各サーバコンピュータ11～14は、図2に示すように構成されている。図2では、サーバコンピュータ11～14の代表として、サーバコンピュータ11の構成を説明する。図2において、サーバコンピュータ11には、主メモリ111、クラスタ共有メモリ用ロック/クラスタ共有ファイルシステム用ロック変換テーブル113、クラスタ共有メモリ用ロック操作手段114、クラスタ共有メモリ用ページ操作手段115、クラスタ共有ファイルシステム用ロック機構116、クラスタ共有ファイルシステム117、クラスタ共有メモリ/クラスタ共有ファイル変換テーブル118、更新ページリスト119、クラスタ共有メモリマップ手段120、クラスタ共有メモリアンマップ手段121、クラスタ共有メモリ用ロックアロケーション手段122、クラスタ共有メモリ用ロック解放手段123が設けられて構成されている。サーバコンピュータ11では、あるアプリケーションプログラムが実行されることで生成されるプロセス112が主メモリ111上のアドレス空間内に配置されて図示しないプロセッサにより実行されて動作している。

10

【0019】

図3は、プロセス112の詳細を示す図である。図3において参照符号Aは、プロセス112のアドレス空間内におけるアドレスを示す。プロセス112には、アドレスAを先頭として複数のページP0、P1、P2、P3、P4、P5が設けられてる。このプロセス112内に仮想的にクラスタ共有メモリが設けられる。

【0020】

クラスタ共有メモリ用ロック/クラスタ共有ファイルシステム用ロック変換テーブル113は、図4に示すようにクラスタ共有メモリ用ロックID欄113aとクラスタ共有ファイルシステム用ロックID欄113bとから構成されている。このクラスタ共有メモリ用ロック/クラスタ共有ファイルシステム用ロック変換テーブル113にクラスタ共有メモリ用ロックIDとクラスタ共有ファイルシステム用ロックIDとを登録することにより、両者の対応を関係づける。

20

【0021】

クラスタ共有メモリ用ロック操作手段114は、クラスタ共有メモリ内に設けられる各ページに記憶されたデータの一貫性を確保するための排他制御をするためにロックを設定（ロックする）し、また設定したロックを解除（アンロック）するための操作手段である。クラスタ共有メモリ用ページ操作手段115は、クラスタ共有メモリに設けられ各ページに対するデータの記憶、各ページに記憶されたデータの読み出しを行う操作手段である。

30

【0022】

クラスタ共有ファイルシステム117は、共有ディスク装置16に記録されたファイルをサーバコンピュータ11、サーバコンピュータ12、サーバコンピュータ13、サーバコンピュータ14とで共有するためのファイルシステムである。クラスタ共有ファイルシステム用ロック機構116は、共有ディスク装置16に記録されたファイルをサーバコンピュータ11、サーバコンピュータ12、サーバコンピュータ13、サーバコンピュータ14とで共有する際に、データの一貫性を確保するための排他制御をするためにロックを設定（ロックする）し、また設定したロックを解除（アンロック）するための操作を行う機構である。

40

【0023】

クラスタ共有メモリ/クラスタ共有ファイル変換テーブル118は、図5に示すように、アドレス欄118a、サイズ欄118b、ファイル名欄118c、ファイル記述子118d、オフセット欄118eとから構成されている。このクラスタ共有メモリ/クラスタ共有ファイル変換テーブル118は、共有ディスク装置16に記録されたファイルをクラスタ共有メモリに対応付け（マッピング）するために設けられている。クラスタ共有メモリマップ手段120は、プロセス112のアドレス空間内にクラスタ共有メモリを仮想的に設けるための手段である。

【0024】

50

クラスタ共有メモリアンマップ手段 1 2 1 は、プロセス 1 1 2 のアドレス空間内に仮想的に設けられたクラスタ共有メモリを解放するための手段である。

クラスタ共有メモリ用ロックアロケーション手段 1 2 2 は、クラスタ共有メモリ用ロックと一対一に対応するクラスタ共有ファイルシステム用ロックとをアロケーションするための手段である。クラスタ共有メモリ用ロック解放手段 1 2 3 は、クラスタ共有メモリ用ロックアロケーション手段 1 2 2 により、アロケーションされたクラスタ共有メモリ用ロックとクラスタ共有ファイルシステム用ロックとを解放するための手段である。

【 0 0 2 5 】

図 6 は、更新ページリスト 1 1 9 の詳細を示す図である。この更新ページリスト 1 1 9 は、共有メモリの中のデータが更新されたページのページ番号を記録するためのリストである。

10

【 0 0 2 6 】

以下、図 7 を用いてサーバーコンピュータ 1 1 で実行されているプロセス 1 1 2 が共有ディスク装置 1 6 に記録されているファイルに対してアクセスする場合の動作を詳細に説明する。

【 0 0 2 7 】

図 7 は、クラスタ共有メモリを用いたクラスタ共有ファイルへのアクセス動作の処理手順を示すフローチャート図である。まず、プロセス 1 1 2 が共有ディスク装置 1 6 に記録されている共有ファイルに対してアクセスしようとする場合には、プロセス 1 1 2 のアドレス空間内に仮想的に共有メモリを設け、この共有メモリ上にアクセス対象の共有ファイルをマッピングする (ステップ S 7 0)。具体的には、プロセス 1 1 2 がクラスタ共有メモリマップ手段 1 2 0 にクラスタ共有メモリマップ操作を指示する。このプロセス 1 1 2 からの指示に基づいてクラスタ共有メモリマップ手段 1 2 0 は、プロセス 1 2 0 が配置された主メモリ 1 1 1 上のアドレス空間内に共有メモリをマッピングする。

20

【 0 0 2 8 】

このクラスタ共有メモリマップ手段 1 2 0 による共有メモリのマッピング操作手順を図 8 に示すフローチャートを用いて説明する。まず、ステップ S 8 0 で、プロセス 1 1 2 のアドレス空間内にクラスタ共有メモリのための領域をアロケーションし、その領域内の全てのページをアクセス不許可に設定する。アロケーションは、プロセス 1 1 2 に C 言語の関数 `m a l l o c ()` を記述することで実施できる。

30

【 0 0 2 9 】

次に、ステップ S 8 1 で、アクセス対象の共有ファイルのプロセス 1 1 2 から指定されたオフセット位置から指定されたサイズ分のデータをステップ S 8 0 でアロケーションした領域にマッピングしたことをクラスタ共有メモリ/クラスタ共有ファイル変換テーブル 1 1 8 に登録する。この登録の結果として、図 5 に示すようにプロセス 1 1 2 のアドレス A からサイズ L の範囲にファイル名が「 D D D D 」で、ファイル記述子が「 7 」で、オフセット「 0 」を登録する。ここまでの処理で、プロセス 1 1 2 が共有ファイルへのアクセスを共有メモリへのアクセスとして処理するための環境を整えたことになる。

【 0 0 3 0 】

続いてステップ S 7 1 において、プロセス 1 1 2 が共有ファイルへのアクセスを共有メモリへのアクセスとして処理するのに必要な排他制御処理をするためにロックを取得する。そこでロック取得の前準備として、プロセス 1 1 2 は、クラスタ共有メモリ用ロックとクラスタ共有ファイル用ロックとを対応づけるために、クラスタ共有メモリ用ロックアロケーション手段 1 2 2 にクラスタ共有メモリ用ロックアロケーション操作を指示する。

40

【 0 0 3 1 】

図 1 0 にクラスタ共有メモリ用ロックアロケーション手段 1 2 2 によるクラスタ共有メモリ用ロックアロケーション操作手順のフローチャートを示す。ステップ S 1 0 0 において、クラスタ共有メモリ用ロックと一対一に対応するクラスタ共有ファイルシステム用ロックをアロケーションするために、図 4 に示すクラスタ共有メモリ用ロック/クラスタ共有ファイルシステム用ロック変換テーブル 1 1 3 にクラスタ共有メモリ用ロックの ID 番号

50

とクラスタ共有ファイルシステム用ロックのID番号とを対応づけて登録する。

【0032】

次にプロセス112は、ロックを取得するために、クラスタ共有メモリ用ロック操作手段114にロックの取得を指示する。このロック取得の指示を受けたクラスタ共有メモリ用ロック操作手段114のロック操作の処理手順を図12に示したフローチャートを用いて説明する。

【0033】

まずステップS120において、クラスタ共有メモリ上の全てのページをアクセス不許可に設定する。次にステップS121において、図4に示したクラスタ共有メモリ用ロック/クラスタファイルシステム用ロック変換テーブル113を用いて、ロック操作されるク
10
ラスタ共有メモリ用ロックをクラスタ共有ファイルシステム用ロックに変換する。続いて、ステップS122において、クラスタ共有メモリ用ロック操作手段114は、クラスタ共有ファイルシステム用ロック機構116に指示して、変換したクラスタファイルシステム用ロックをロック操作してロックを取得する。

【0034】

ステップS71でロック取得が成功した場合には、ステップS72に進む。
もし、ロック取得が失敗した場合には、換言すると既に他のプロセスが同様な処理により、実体としての共有ファイルのロックを取得している場合には、ロックが取得できないので、ロックが解放されて、取得できるまで処理を待機しステップS71の処理を再び続ける
20

【0035】

続いてステップS72において、プロセス112は、クラスタ共有メモリ用ページ操作手段115に指示して、共有メモリへのREADアクセス（共有メモリに対するload命令の実行）又はWRITEアクセス（共有メモリに対するstore命令の実行）を実行させる。このクラスタ共有メモリ用ページ操作手段115が共有メモリへアクセスを実行すると、この共有メモリの全てのページはアクセス不許可に設定されているため、ページフォールトが発生する。クラスタ共有メモリ用ページ操作手段115がREADアクセスした場合には、READページフォールトが発生する。また、クラスタ共有メモリ用ページ操作手段115がWRITEアクセスした場合には、WRITEページフォールトが発生する。
30

【0036】

READページフォールトが発生した場合のクラスタ共有メモリ用ページ操作手段115の処理手順を図14に示す。WRITEページフォールトが発生した場合のクラスタ共有メモリ用ページ操作手段115の処理手順を図15に示す。

【0037】

図14を参照して、READページフォールトが発生した場合のクラスタ共有メモリ用ページ操作手段115の処理手順を説明する。まず、ステップS140で、ページフォールトが発生した共有メモリ内のページのデータをクラスタ共有ファイルシステム117によりクラスタ共有ファイル（共有ディスク装置16）の対応する部分から当該共有メモリ内のページに読み込む。この時、ページのデータをクラスタ共有ファイルのどの部分から読み込むかは、クラスタ共有メモリ/クラスタ共有ファイル変換テーブル118に登録されたデータに基づいて求める。続いて、ステップS141で、ページフォールトが発生したページをREAD可能に設定する。
40

【0038】

次に図15を参照して、WRITEページフォールトが発生した場合のクラスタ共有メモリ用ページ操作手段115の処理手順を説明する。まず、ステップS150において、ページフォールトが発生したページがREAD可能に設定されているかどうかを調べる。READ可能な場合には、ステップS152へ進む。READ可能でない場合には、ステップS151に進む。

【0039】

10

20

30

40

50

ステップS 1 5 1では、ページフォールトが発生した共有メモリ内のページのデータをクラスタ共有ファイルシステム 1 1 7によりクラスタ共有ファイル(共有ディスク装置 1 6)の対応する部分から当該共有メモリ内のページに読み込む。この時、ページのデータをクラスタ共有ファイルのどの部分から読み込むかは、クラスタ共有メモリ/クラスタ共有ファイル変換テーブル 1 1 8に登録されたデータに基づいて求める。

【 0 0 4 0 】

続いて、ステップS 1 5 2で、ページフォールトが発生した共有メモリ内のページをREAD/WRITE可能に設定する。最後に、ステップS 1 5 3で、ページフォールトが発生した共有メモリ内のページの番号を図 6 に示す更新ページリスト 1 1 9に登録する。クラスタ共有メモリ用ページ操作手段 1 1 5は、このようにページフォールト処理をした後、実際にクラスタ共有メモリへアクセスする。

10

【 0 0 4 1 】

このようにステップS 7 2で共有メモリへのアクセスをした後、続くステップS 7 3において、プロセス 1 1 2は、クラスタ共有メモリ用ロックを解放するためにクラスタ共有メモリ用ロック操作手段 1 1 4に指示して、クラスタ共有メモリロックのアンロック操作を実行させる。

【 0 0 4 2 】

ここで図 1 3を参照して、クラスタ共有メモリ用ロック操作手段 1 1 4によるクラスタ共有メモリロックのアンロック操作手順を説明する。図 1 3は、クラスタ共有メモリ用ロック操作手段 1 1 4によるクラスタ共有メモリロックのアンロック操作の操作手順を示すフローチャート図である。

20

【 0 0 4 3 】

まず、ステップS 1 3 0において、図 6 に示した更新ページリスト 1 1 9に登録されている全ての更新ページのデータをクラスタ共有ファイルシステム 1 1 7により共有ディスク装置 1 6に記録されているクラスタ共有ファイルの該当部分に書き込む。この時、更新ページをクラスタ共有ファイルのどの部分に書き込むかは、クラスタ共有メモリ/クラスタ共有ファイル変換テーブル 1 1 8を用いて求める。続いて、ステップS 1 3 1において、更新ページリスト 1 1 9をクリアする。

【 0 0 4 4 】

以上で図 7 に示したクラスタ共有メモリを用いたクラスタ共有ファイルへのアクセス動作の処理手順の説明を終了する。ここで、プロセス 1 1 2は、これまでの説明でアクセスした共有ファイルをまだアクセスする場合には、図 4 に示したクラスタ共有メモリ用ロック/クラスタ共有ファイル用ロック変換テーブル 1 1 3のエントリを削除せずに残しておく。このエントリを残しておけば、次にロックを取得する場合には、クラスタ共有メモリ用ロック操作手段 1 1 4がクラスタ共有ファイルシステム用ロック機構 1 1 6にファイルのロック取得を依頼するだけでロック取得ができる。

30

【 0 0 4 5 】

もし、プロセス 1 1 2がこれまでの説明でアクセスした共有ファイルを以後アクセスすることがない場合には、プロセス 1 1 2は、クラスタ共有メモリ用ロック解放手段 1 2 3に指示してクラスタ共有メモリ用ロックを解放させる。図 1 1 は、クラスタ共有メモリ用ロック解放手段 1 2 3によるクラスタ共有メモリ用ロックを解放させる処理手順を示すフローチャート図である。図 1 1 において、ステップS 1 1 0で、プロセス 1 1 2が指定したクラスタ共有メモリ用ロックに関するエントリを図 4 に示すクラスタ共有メモリ用ロック/クラスタ共有ファイルシステム用ロック変換テーブル 1 1 3から抹消する。

40

これにより、クラスタ共有メモリ用ロックが解放される。

【 0 0 4 6 】

最後にプロセス 1 1 2が共有メモリへのアクセスとして処理すべき全ての共有ファイルへのアクセスが終了した場合には、プロセス 1 1 2がクラスタ共有メモリアンマップ手段 1 2 1に指示してプロセスのアドレス空間内に仮想的に設定した共有メモリのマッピングを

50

解除させる。すなわち、クラスタ共有メモリをアンマップさせる。

【 0 0 4 7 】

図9は、クラスタ共有メモリアンマップ手段121がクラスタ共有メモリをアンマップする処理手順を示すフローチャート図である。図9において、まず、ステップS90において、プロセス112のアドレス空間内にアロケーションしたクラスタ共有メモリのための領域を解放する。続いて、ステップS91において、アンマップする領域に関するエントリを図5に示すクラスタ共有メモリ/クラスタ共有ファイル変換テーブル118から抹消する。

【 0 0 4 8 】

次に図16を用いて本発明において、2つのサーバーコンピュータが共有ディスク装置に記録されている同じファイルに対するアクセスをする場合の動作を説明する。図16は、サーバーコンピュータAとサーバーコンピュータBと共有ディスク装置Cとから構成されるクラスタシステムの動作を説明するためのシステム構成を示す図である。

【 0 0 4 9 】

図16において、サーバーコンピュータAとサーバーコンピュータBとは、それぞれ共有ディスク装置Cに接続されている。サーバーコンピュータAでは、プロセスAが動作している。また、サーバーコンピュータBでは、プロセスBが動作している。これらのプロセスAとプロセスBは、それぞれ図示を省略したクラスタ共有ファイルシステムで管理されている共有ディスク装置Cに記録されたクラスタ共有ファイルDを先に説明したように自身のアドレス空間内にマッピングしているものとする。クラスタ共有ファイルDの中にはデータ領域Xが存在する。

【 0 0 5 0 】

このように設定された状況の下に、プロセスAが(1)~(3)の処理を実行し、次にプロセスBが(4)~(6)の処理を実行し、更にプロセスA(7)~(9)の処理を実行するとする。まずプロセスAは(1)でクラスタ共有メモリ用ロックを取得する。これによりクラスタ共有ファイルDがクラスタ共有メモリとしてマップされた領域のページは全てアクセス不許可になる。次に(2)でクラスタ共有メモリ上のデータ領域Xに数値「1」を加えようとする。するとWRITEページフォルトが発生する。クラスタ共有メモリ操作手段は、この領域がREAD不可能なので、このデータ領域を含むページの内容をクラスタ共有ファイルDから読みこむ。

【 0 0 5 1 】

そして、そのページをREAD/WRITE可能に設定し、ページフォルト処理を終了する。ページフォルトから戻った後、プロセスAはデータ領域Xの値に1を加える。そして(3)でクラスタ共有メモリ用ロックを解放する。

これにより、クラスタ共有メモリ用ロック操作手段は、更新ページであるデータ領域Xを含むページの内容をクラスタ共有ファイルDに書き戻す。データ領域Xの初期値が「0」だったとすると、この時点でデータ領域Xの値は「1」となる。

【 0 0 5 2 】

次にプロセスBは(4)でクラスタ共有メモリ用ロックを取得する。

これによりクラスタ共有ファイルD4がクラスタ共有メモリとしてマップされた領域のページは全てアクセス不許可になる。次に(5)でクラスタ共有メモリ上のデータ領域Xに数値「1」を加えようとする。するとWRITEページフォルトが発生する。クラスタ共有メモリ操作手段は、この領域がREADが不可能なので、このデータ領域を含むページの内容をクラスタ共有ファイルDから読みこむ。そして、そのページをREAD/WRITE可能に設定し、ページフォルト処理を終了する。ページフォルトから戻った後、プロセスBはデータ領域Xの値に数値「1」を加える。そして(6)でクラスタ共有メモリ用ロックを解放する。これにより、クラスタ共有メモリ用ロック操作手段は、更新ページであるデータ領域Xを含むページの内容をクラスタ共有ファイルDに書き戻す。この時点でデータ領域Xの値は「2」になる。

【 0 0 5 3 】

最後にプロセスAは(7)でクラスタ共有メモリ用ロックを取得する。これによりクラスタ共有ファイルDがクラスタ共有メモリとしてマップされた領域のページは全てアクセス不許可になる。次に(8)でクラスタ共有メモリ上のデータ領域Xに1を加えようとする。するとWRITEページフォールトが発生する。クラスタ共有メモリ操作手段は、この領域がREAD/WRITE不可能なので、このデータ領域を含むページの内容をクラスタ共有ファイルDから読みこむ。そして、そのページをすREAD/WRITE可能に設定し、ページフォールト処理を終了する。ページフォールトから戻った後、プロセスAはデータ領域Xの値に数値「1」を加える。そして(9)でクラスタ共有メモリ用ロックを解放する。これにより、クラスタ共有メモリ用ロック操作手段は、更新ページであるデータ領域Xを含むページの内容をクラスタ共有ファイルDに書き戻す。この時点でデータ領域Xの値は3になる。

10

【0054】

【発明の効果】

本発明を適用することにより、クラスタ共有ファイルシステムを持つクラスタシステムで、ファイルへの共有アクセスと共に、メモリへの共有アクセスも可能になる。

更に本発明では、クラスタ共有ファイルシステムを用いることで、クラスタ共有メモリを安価に実現することができる。また、クラスタ共有ファイルをクラスタ共有メモリとしてマッピングするので、更新データに永続性が持たれる。更に、本発明では、クラスタ共有メモリを、クラスタ共有ファイルシステムを用いて実現するため、クラスタ共有ファイルを分散共有メモリとしてマッピングし、メモリとしてのアクセス(load命令/sto
re命令)とファイルとしてのアクセス(readシステムコール/writeシステム
コール)を並列に実行することができるようになる。

20

【図面の簡単な説明】

【図1】本発明のクラスタシステム10を示す図である。

【図2】サーバーコンピュータ11の構成を示す図である。

【図3】プロセス112の詳細を示した図である。

【図4】クラスタ共有メモリ用ロック/クラスタ共有ファイルシステム用ロック変換テーブル113の詳細を示した図である。

【図5】クラスタ共有メモリ/クラスタ共有ファイル変換テーブル118の詳細を示した図である。

30

【図6】更新ページリスト119の詳細を示した図である。

【図7】クラスタ共有メモリを用いたクラスタ共有ファイルへのアクセス動作の処理手順を示すフローチャート図である。

【図8】クラスタ共有メモリマップ手段120による共有メモリのマッピング操作手順を示すフローチャート図である。

【図9】クラスタ共有メモリアンマップ手段121がクラスタ共有メモリをアンマップする処理手順を示すフローチャート図である。

【図10】クラスタ共有メモリ用ロックアロケーション手段122によるクラスタ共有メモリ用ロックアロケーション操作手順を示すフローチャート図である。

【図11】クラスタ共有メモリ用ロック解放手段123によるクラスタ共有メモリ用ロックを解放させる処理手順を示すフローチャート図である。

40

【図12】クラスタ共有メモリ用ロック操作手段114のロック操作の処理手順を示すフローチャート図である。

【図13】クラスタ共有メモリ用ロック操作手段114によるクラスタ共有メモリロックのアンロック操作の処理手順を示すフローチャート図である。

【図14】READページフォールトが発生が発生した場合のクラスタ共有メモリ用ページ操作手段115の処理手順を示すフローチャート図である。

【図15】WRITEページフォールトが発生が発生した場合のクラスタ共有メモリ用ページ操作手段115の処理手順を示すフローチャート図である。

【図16】2つのサーバーコンピュータが共有ディスク装置に記録されている同じファイ

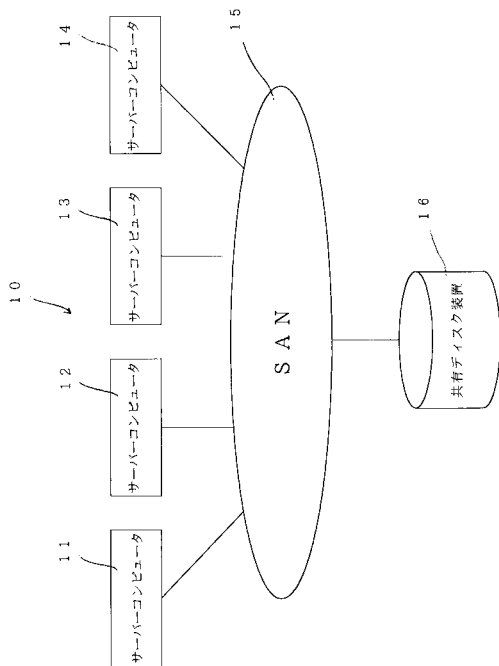
50

ルに対するアクセスをする場合の動作を説明するためのシステム構成を示す図である。

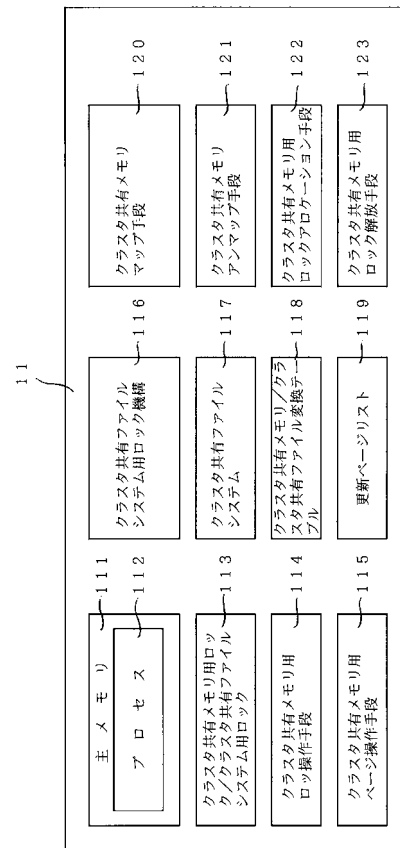
【符号の説明】

- 1 1 サーバコンピュータ
- 1 2 サーバコンピュータ
- 1 3 サーバコンピュータ
- 1 4 サーバコンピュータ
- 1 5 S A N
- 1 6 共有ディスク装置

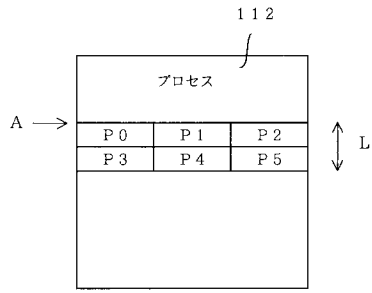
【図 1】



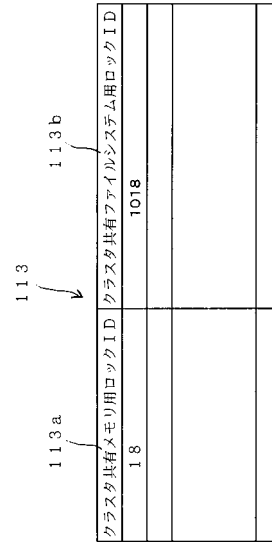
【図 2】



【図3】



【図4】

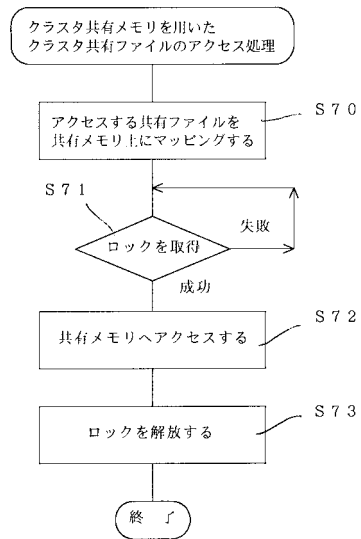


【図5】

118

118a	118b	118c	118d	118e
アドレス	サイズ	ファイル名	ファイル記述	オフセット
A	L	DDDD	7	0

【図7】

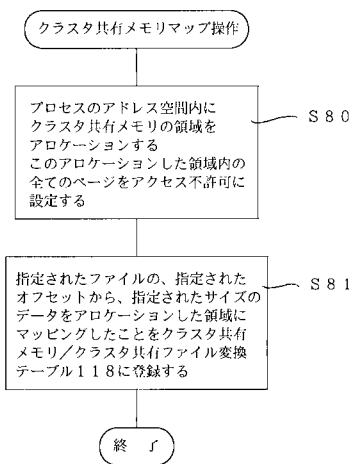


【図6】

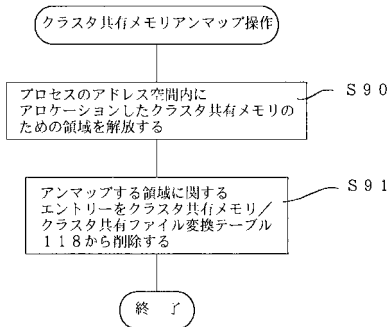
119

P4	P2		

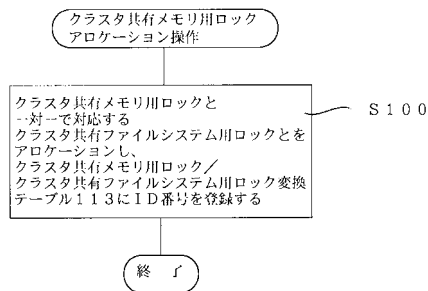
【 図 8 】



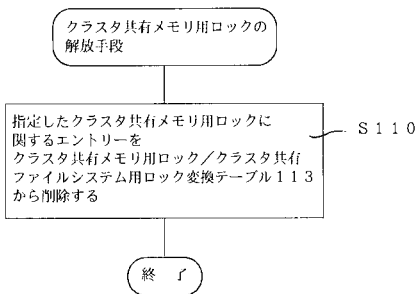
【 図 9 】



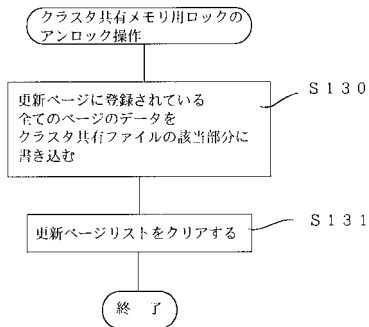
【 図 10 】



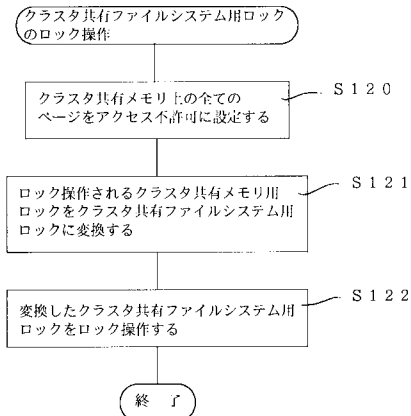
【 図 11 】



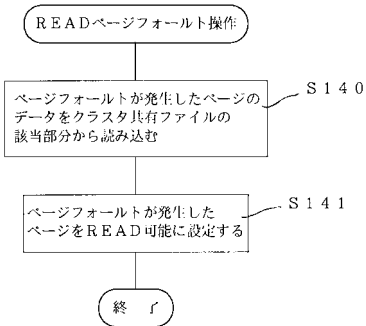
【 図 13 】



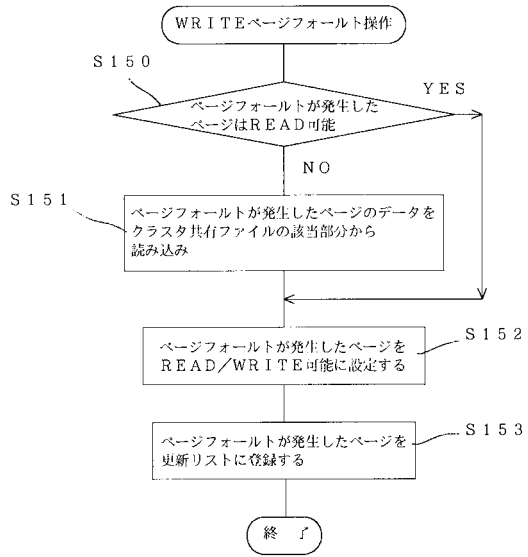
【 図 12 】



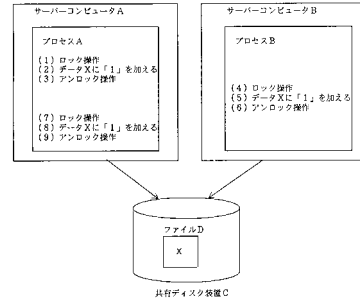
【 図 14 】



【 図 15 】



【 図 16 】



フロントページの続き

(74)代理人 100084618

弁理士 村松 貞男

(74)代理人 100092196

弁理士 橋本 良郎

(72)発明者 平山 秀昭

東京都府中市東芝町1番地 株式会社東芝 府中事業所内

審査官 平井 誠

(56)参考文献 特開2000-305832(JP,A)

米国特許出願公開第2002/0133675(US,A1)

Shared Memory RAM Disk for a Cluster with Shared Memory, IBM Technical Disclosure Bulletin, 米国, 1993年, Vol.36, No.06B, p.299-300

(58)調査した分野(Int.Cl.⁷, DB名)

G06F 12/00