

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3815967号

(P3815967)

(45) 発行日 平成18年8月30日(2006.8.30)

(24) 登録日 平成18年6月16日(2006.6.16)

(51) Int. Cl.

F I

G 0 6 F 12/00 (2006.01)

G 0 6 F 12/08 (2006.01)

G O 6 F 12/00 5 3 5 Z

G O 6 F 12/00 5 1 4 M

G O 6 F 12/00 5 4 5 A

G O 6 F 12/08 5 3 1 Z

G O 6 F 12/08 5 5 7

請求項の数 190 (全 42 頁)

(21) 出願番号 特願2000-531781 (P2000-531781)
 (86) (22) 出願日 平成11年2月12日(1999.2.12)
 (65) 公表番号 特表2002-503846 (P2002-503846A)
 (43) 公表日 平成14年2月5日(2002.2.5)
 (86) 国際出願番号 PCT/US1999/002965
 (87) 国際公開番号 W01999/041664
 (87) 国際公開日 平成11年8月19日(1999.8.19)
 審査請求日 平成16年7月28日(2004.7.28)
 (31) 優先権主張番号 60/074,587
 (32) 優先日 平成10年2月13日(1998.2.13)
 (33) 優先権主張国 米国(US)
 (31) 優先権主張番号 09/199,120
 (32) 優先日 平成10年11月24日(1998.11.24)
 (33) 優先権主張国 米国(US)

(73) 特許権者 502303739
 オラクル・インターナショナル・コーポレ
 イション
 アメリカ合衆国、94065 カリフォル
 ニア州、レッドウッド・ショアーズ、オラ
 クル・パークウェイ、500
 (74) 代理人 100064746
 弁理士 深見 久郎
 (74) 代理人 100085132
 弁理士 森田 俊雄
 (74) 代理人 100083703
 弁理士 仲村 義平
 (74) 代理人 100096781
 弁理士 堀井 豊

最終頁に続く

(54) 【発明の名称】 あるノードのキャッシュから別のノードのキャッシュヘデータを転送するための方法および装置

(57) 【特許請求の範囲】

【請求項1】

第1のキャッシュから第2のキャッシュへ資源を転送するための方法であって、
 前記第1のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することな
 しに、前記第1のキャッシュ内に資源の第1のコピーを保持する一方で資源の第2のコピ
 ーを第1のキャッシュから第2のキャッシュへと転送するステップと、
 資源の前記第1のコピーまたはそのサクセサが永続的に記憶されるまで前記第1のキャ
 ッシュ内に資源の少なくとも1つのコピーを保持するステップとを含む、方法。

【請求項2】

前記第1のキャッシュは第1のデータベースサーバによって維持されるキャッシュであり、
 前記第2のキャッシュは第2のデータベースサーバによって維持されるキャッシュであ
 る、請求項1に記載の方法。

【請求項3】

前記第2のコピーを前記第2のキャッシュに転送するよりも前に前記資源の前記第1の
 コピーが前記第1のキャッシュ内で変更されることを可能にするステップと、

前記第2のコピーを前記第2のキャッシュに転送した後に前記資源の前記第1のコピー
 が変更されることを防ぐステップとをさらに含む、請求項1に記載の方法。

【請求項4】

前記第2のコピーを前記第2のキャッシュに転送した後に、前記第1のコピーの開放の
 許可の要求を送信するステップと、

10

20

前記要求に応答して、前記第 1 のコピーまたはそのサクセサが永続的に記憶されるようにするステップと、

前記サクセサが永続的に記憶されたことに応答して、前記第 1 のコピーが開放可能であることを示すメッセージを送信するステップとをさらに含む、請求項 1 に記載の方法。

【請求項 5】

前記第 1 のコピーの開放の許可の要求を送信するステップは、送信側プロセスによって実行され、

前記第 1 のコピーまたはそのサクセサが永続的に記憶されるようにするステップは、送信側プロセス以外のプロセスが前記資源の前記第 1 のコピーのサクセサを記憶するようにするステップを含む、請求項 4 に記載の方法。

10

【請求項 6】

前記第 1 のキャッシュ内に資源の少なくとも 1 つのコピーを保持するステップは、

前記第 1 のコピーを永続的に記憶しようとする前に、前記資源の永続的に記憶されたコピーが前記第 1 のコピーより最近のものであるかどうか決定するステップと、

もし前記永続的に記憶されたコピーが前記第 1 のコピーより最近のものであれば、前記第 1 のコピーを永続的に記憶することなしに前記第 1 のコピーを開放するステップと、

もし前記永続的に記憶されたコピーが前記第 1 のコピーより最近のもでなければ、前記第 1 のコピーを永続的に記憶するステップとを含む、請求項 1 に記載の方法。

【請求項 7】

変更許可を、第 1 のキャッシュに関連付けられる送信側プロセスから第 2 のキャッシュに関連付けられる受信側プロセスに、前記資源の前記第 2 のコピーとともに転送するステップをさらに含む、請求項 3 に記載の方法。

20

【請求項 8】

前記資源にアクセスする許可はマスタによって管理され、

前記変更許可を前記受信側プロセスに転送するステップは、前記変更許可の前記受信側プロセスへの転送についての確認を前記マスタから受信するより前に実行される、請求項 7 に記載の方法。

【請求項 9】

変更許可を、前記第 1 のキャッシュに関連付けられる送信側プロセスから前記第 2 のキャッシュに関連付けられる受信側プロセスに、前記資源の第 2 のコピーとともに転送するステップをさらに含み、

30

前記資源にアクセスする許可はマスタによって管理され、

前記変更許可を前記受信側プロセスに転送するステップは、前記変更許可の前記受信側プロセスへの転送の確認を前記マスタが受信するより前に実行される、請求項 1 に記載の方法。

【請求項 10】

前記第 2 のキャッシュに関連付けられる受信側プロセスが前記資源の要求を前記資源のマスタに送信するステップと、

前記受信側プロセスからの前記要求に応答して、前記資源の前記マスタが前記第 1 のキャッシュに関連付けられる送信側プロセスにメッセージを送信するステップと、

40

前記送信側プロセスが前記マスタからの前記メッセージに応答して、前記第 2 のコピーを前記受信側プロセスに転送するステップとをさらに含む、請求項 1 に記載の方法。

【請求項 11】

前記第 2 のコピーを前記第 2 のキャッシュに転送するステップの後に、

前記第 1 のキャッシュに関連付けられる送信側プロセスが、ロックマネージャにロックを要求するステップを含み、前記ロックは前記資源をディスクに書込む許可は与えるが前記資源を変更する許可は与えず、さらに、

前記ロックマネージャが、前記第 1 のコピーと少なくとも同じほど最近である前記資源のバージョンを有するプロセスを選択するステップと、

前記ロックマネージャが、前記ロックを前記選択されたプロセスに与えるステップと、

50

前記選択されたプロセスが、前記資源の前記バージョンをディスクに書込むステップとをさらに含む、請求項 1 に記載の方法。

【請求項 1 2】

前記資源のバージョンがディスクに書込まれるのに応答して、前記ロックマネージャが前記バージョンよりも古い前記資源のすべてのバージョンが開放されるようにするステップをさらに含む、請求項 1 1 に記載の方法。

【請求項 1 3】

前記資源のダーティコピーを保持するキャッシュの障害の後に、
前記障害の発生したキャッシュが資源の最新バージョンを保持していたかどうかを決定するステップと、

10

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していたならば、
資源の最新パストイメージをディスクに書込むステップと、
資源の以前のすべてのパストイメージを開放するステップと、
前記障害の発生したキャッシュの復旧ログを適用して資源の最新バージョンを再構築するステップとをさらに含む、請求項 1 に記載の方法。

【請求項 1 4】

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していなかったならば、

資源の最新バージョンをディスクに書込むステップと、
資源のすべてのパストイメージを開放するステップとをさらに含む、請求項 1 3 に記載の方法。

20

【請求項 1 5】

前記資源のダーティバージョンを保持する複数のキャッシュの障害の後に、
前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

前記障害の発生したキャッシュの復旧ログをマージし適用して資源の最新バージョンを再構築するステップをさらに含む、請求項 1 に記載の方法。

【請求項 1 6】

30

複数のキャッシュが 1 つ以上の共有されるディスクからの資源のダーティバージョンを保持するシステムにおいてデータを管理する方法であって、

前記複数のキャッシュのうち 1 つのキャッシュが、前記複数のキャッシュのうち別のキャッシュにある資源のダーティバージョンを要求するとき、前記資源のダーティバージョンが存在するキャッシュから、前記ダーティバージョンを、前記ダーティバージョンを要求するキャッシュに、前記ダーティバージョンを最初に永続的に記憶することなしに転送するステップと、

前記複数のキャッシュの各キャッシュに対して別々の復旧ログを維持するステップと、
前記複数のキャッシュのうち 1 つのキャッシュに障害が発生したとき、前記障害の発生したキャッシュを、前記障害の発生したキャッシュに関連付けられた復旧ログに基づいて、
前記複数のキャッシュのうちの他のキャッシュの前記別々の復旧ログを検査することなしに、復旧させるステップとを含む、方法。

40

【請求項 1 7】

前記複数のキャッシュの各キャッシュは、複数のデータベースサーバの別々のデータベースサーバによって維持されるキャッシュである、請求項 1 6 に記載の方法。

【請求項 1 8】

前記ダーティバージョンが存在するキャッシュは第 1 のキャッシュであり、
前記第 1 のキャッシュ内の資源のダーティバージョンは資源の第 1 のコピーであり、
前記ダーティバージョンを要求するキャッシュは第 2 のキャッシュであり、
前記ダーティバージョンを転送するステップは、前記資源の第 2 のコピーを前記第 2 の

50

キャッシュに転送することにより実行され、

前記方法はさらに、

前記第 2 のコピーを前記第 2 のキャッシュに転送するよりも前に前記資源の前記第 1 のコピーが前記第 1 のキャッシュ内で変更されることを可能にするステップと、

前記第 2 のコピーを前記第 2 のキャッシュに転送した後に前記資源の前記第 1 のコピーが変更されることを防ぐステップとをさらに含む、請求項 16 に記載の方法。

【請求項 19】

複数のノードによって共有される資源を管理するための方法であって、

前記複数のノードのうち第 1 のノードの第 1 のキャッシュ内で前記資源を変更して前記資源の変更されたバージョンを作成するステップと、

前記第 1 のノードに障害が発生したときにどこで作業を開始するかを示す、前記第 1 のノードについてのチェックポイントを維持するステップと、

前記第 1 のキャッシュから永続的記憶装置に前記変更されたバージョンを最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に前記変更されたバージョンの第 1 のコピーを保持する一方で、前記変更されたバージョンの第 2 のコピーを前記第 1 のキャッシュから前記複数のノードのうち第 2 のノードの第 2 のキャッシュに転送するステップと、

前記複数のノードのうち別のノードが、前記変更されたバージョンと少なくとも同じほど最近のものである前記資源のバージョンを永続的に記憶したことが示されることに応答して、前記チェックポイントを進めるステップとを含む、方法。

【請求項 20】

前記変更されたバージョンの前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで、前記第 1 のキャッシュ内に前記変更されたバージョンの少なくとも 1 つのコピーを保持するステップをさらに含む、請求項 19 に記載の方法。

【請求項 21】

前記複数のノードのうち別のノードが、前記変更されたバージョンと少なくとも同じほど最近のものである前記資源のバージョンを永続的に記憶したことが示されることに応答して、前記第 1 のノードで前記資源を開放するステップを含む、請求項 19 に記載の方法。

【請求項 22】

第 1 のキャッシュから第 2 のキャッシュへ資源を転送するための方法であって、

前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから前記第 2 のキャッシュへと転送するステップと、

前記資源の第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 1 のコピーが前記第 1 のキャッシュ内で置換されることを防ぐステップとを含む、方法。

【請求項 23】

命令の 1 つ以上のシーケンスを搬送して第 1 のキャッシュから第 2 のキャッシュへ資源を転送するためのコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサは、

前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを第 1 のキャッシュから第 2 のキャッシュへと転送するステップと、

資源の前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 1 のキャッシュ内に資源の少なくとも 1 つのコピーを保持するステップとを実行する、コンピュータ読出可能媒体。

【請求項 24】

前記第 2 のコピーを前記第 2 のキャッシュに転送するよりも前に前記資源の前記第 1 のコピーが前記第 1 のキャッシュ内で変更されることを可能にするステップと、

前記第 2 のコピーを前記第 2 のキャッシュに転送した後に前記資源の前記第 1 のコピーが変更されることを防ぐステップとを実行するための命令のシーケンスをさらに含む、請

10

20

30

40

50

求項 2 3 に記載のコンピュータ読出可能媒体。

【請求項 2 5】

前記第 2 のコピーを前記第 2 のキャッシュに転送した後に前記第 1 のコピーの開放の許可の要求を送信するステップと、

前記要求に応答して、前記第 1 のコピーまたはそのサクセサが永続的に記憶されるようにするステップと、

前記サクセサが永続的に記憶されたことに応答して、前記第 1 のコピーが置換可能であることを示すメッセージを送信するステップとを実行するための命令のシーケンスをさらに含む、請求項 2 3 に記載のコンピュータ読出可能媒体。

【請求項 2 6】

前記第 1 のコピーの開放の許可の要求を送信するステップは送信側プロセスによって実行され、

前記第 1 のコピーまたはそのサクセサが永続的に記憶されるようにするステップは、前記送信側プロセス以外のプロセスが前記資源の前記第 1 のコピーのサクセサを記憶するようにするステップを含む、請求項 2 5 に記載のコンピュータ読出可能媒体。

【請求項 2 7】

前記第 1 のキャッシュ内に資源の少なくとも 1 つのコピーを保持するステップは、

前記第 1 のコピーを永続的に記憶しようとする前に、前記資源の永続的に記憶されたコピーは前記第 1 のコピーより最近のものであるかどうか決定するステップと、

もし前記永続的に記憶されたコピーが前記第 1 のコピーより最近のものであれば、前記第 1 のコピーを永続的に記憶することなしに、前記第 1 のコピーを開放するステップと、

もし前記永続的に記憶されたコピーが前記第 1 のコピーよりも最近のもでなければ、前記第 1 のコピーを永続的に記憶するステップとを含む、請求項 2 3 に記載のコンピュータ読出可能媒体。

【請求項 2 8】

変更許可を、前記第 1 のキャッシュに関連付けられる送信側プロセスから前記第 2 のキャッシュに関連付けられる受信側プロセスに、前記資源の前記第 2 のコピーとともに転送するステップを実行するための命令のシーケンスをさらに含む、請求項 2 4 に記載のコンピュータ読出可能媒体。

【請求項 2 9】

前記資源にアクセスする許可はマスタによって管理され、

前記変更許可を受信側プロセスに転送するステップは、前記変更許可の前記受信側プロセスへの転送についての確認を前記マスタから受信するより前に実行される、請求項 2 8 に記載のコンピュータ読出可能媒体。

【請求項 3 0】

変更許可を、前記第 1 のキャッシュに関連付けられる送信側プロセスから前記第 2 のキャッシュに関連付けられる受信側プロセスに、前記資源の第 2 のコピーとともに転送するステップを実行するための命令のシーケンスをさらに含む、

前記資源にアクセスする許可はマスタによって管理され、

前記変更許可を前記受信側プロセスに転送するステップは、前記変更許可の前記受信側プロセスへの転送の確認を前記マスタが受信するより前に実行される、請求項 2 3 に記載のコンピュータ読出可能媒体。

【請求項 3 1】

前記第 2 のキャッシュに関連付けられる受信側プロセスが前記資源の要求を前記資源のマスタに送信するステップと、

前記受信側プロセスからの前記要求に応答して、前記資源の前記マスタが前記第 1 のキャッシュに関連付けられる送信側プロセスにメッセージを送信するステップと、

前記送信側プロセスが前記マスタからの前記メッセージに응答して、前記第 2 のコピーを前記受信側プロセスに転送するステップとを実行するための命令のシーケンスをさらに含む、請求項 2 3 に記載のコンピュータ読出可能媒体。

10

20

30

40

50

【請求項 3 2】

前記第 2 のコピーを前記第 2 のキャッシュに転送するステップの後に、

前記第 1 のキャッシュに関連付けられる送信側プロセスが、ロックマネージャにロックを要求するステップを実行するための命令のシーケンスをさらに含み、前記ロックは前記資源をディスクに書込む許可は与えるが前記資源を変更する許可は与えず、さらに、

前記ロックマネージャが、前記第 1 のコピーと少なくとも同じほど最近である前記資源のバージョンを有するプロセスを選択するステップと、

前記ロックマネージャが、前記ロックを前記選択されたプロセスに与えるステップと、

前記選択されたプロセスが、前記資源の前記バージョンをディスクに書込むステップとを実行するための命令のシーケンスを含む、請求項 2 3 に記載のコンピュータ読出可能媒体。

10

【請求項 3 3】

前記資源の前記バージョンがディスクに書込まれるのに応答して、前記ロックマネージャが前記バージョンよりも古い前記資源のすべてのバージョンが開放されるようにするステップを実行するための命令のシーケンスをさらに含む、請求項 3 2 に記載のコンピュータ読出可能媒体。

【請求項 3 4】

前記資源のダークティコピーを保持するキャッシュの障害の後に、

前記障害の発生したキャッシュが資源の最新バージョンを保持していたかどうかを決定するステップと、

20

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していたならば、

資源の最新パストイメージをディスクに書込むステップと、

資源の以前のすべてのパストイメージを開放するステップと、

前記障害の発生したキャッシュの復旧ログを適用して資源の最新バージョンを再構築するステップとを実行するための命令のシーケンスをさらに含む、請求項 2 3 に記載のコンピュータ読出可能媒体。

【請求項 3 5】

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していなかったならば、

資源の最新バージョンをディスクに書込むステップと、

30

資源のすべてのパストイメージを開放するステップとを実行するための命令のシーケンスをさらに含む、請求項 3 4 に記載のコンピュータ読出可能媒体。

【請求項 3 6】

前記資源のダークティバージョンを保持する複数のキャッシュの障害の後に、

前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

前記障害の発生したキャッシュの復旧ログをマージし適用して資源の最新バージョンを再構築するステップとを実行するための命令のシーケンスをさらに含む、請求項 2 3 に記載のコンピュータ読出可能媒体。

40

【請求項 3 7】

前記第 1 のキャッシュ内に資源の第 1 のコピーを保持するための命令は、第 1 のデータベースサーバによって維持される第 1 のキャッシュ内に資源の第 1 のコピーを保持するための命令を含み、

前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに資源の第 2 のコピーを第 1 のキャッシュから第 2 のキャッシュへと転送するための命令は、前記第 1 のデータベースサーバによって維持される前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに資源の第 2 のコピーを前記第 1 のデータベースサーバによって維持される第 1 のキャッシュから第 2 のデータベースサ

50

ーバによって維持される第2のキャッシュへと転送するための命令を含み、

資源の前記第1のコピーまたはそのサクセサが永続的に記憶されるまで前記第1のキャッシュ内に資源の少なくとも1つのコピーを保持するための命令は、資源の前記第1のコピーまたはそのサクセサが永続的に記憶されるまで前記第1のデータベースサーバによって維持される前記第1のキャッシュ内に資源の少なくとも1つのコピーを保持するための命令を含む、請求項23に記載のコンピュータ読出可能媒体。

【請求項38】

命令を搬送し、複数のキャッシュが1つ以上の共有されるディスクからの資源のダーティバージョンを保持するシステムにおいてデータを管理するためのコンピュータ読出可能媒体であって、前記命令は、

前記複数のキャッシュのうち1つのキャッシュが、前記複数のキャッシュのうち別のキャッシュにある資源のダーティバージョンを要求するとき、前記資源のダーティバージョンが存在するキャッシュから、前記ダーティバージョンを、前記ダーティバージョンを要求するキャッシュに、前記ダーティバージョンを最初に永続的に記憶することなしに転送するステップと、

前記複数のキャッシュの各キャッシュに対して別々の復旧ログを維持するステップと、
前記複数のキャッシュのうち1つのキャッシュに障害が発生したとき、前記障害の発生したキャッシュを、前記障害の発生したキャッシュに関連付けられた復旧ログに基づいて、前記複数のキャッシュのうちの他のキャッシュの前記別々の復旧ログを検査することなしに、復旧させるステップとを実行するための命令を含む、コンピュータ読出可能媒体。

【請求項39】

前記複数のキャッシュの各キャッシュは、複数のデータベースサーバの別々のデータベースサーバによって維持されるキャッシュである、請求項38に記載のコンピュータ読出可能媒体。

【請求項40】

前記ダーティバージョンが存在するキャッシュは第1のキャッシュであり、
前記第1のキャッシュ内の資源のダーティバージョンは資源の第1のコピーであり、
前記ダーティバージョンを要求するキャッシュは第2のキャッシュであり、
前記ダーティバージョンを転送するステップは、前記資源の第2のコピーを前記第2のキャッシュに転送することにより実行され、

前記コンピュータ読出可能媒体はさらに、
前記第2のコピーを前記第2のキャッシュに転送するよりも前に前記資源の前記第1のコピーが前記第1のキャッシュ内で変更されることを可能にするステップと、

前記第2のコピーを前記第2のキャッシュに転送した後に前記資源の前記第1のコピーが変更されることを防ぐステップとを実行するための命令をさらに含む、請求項38に記載のコンピュータ読出可能媒体。

【請求項41】

命令の1つ以上のシーケンスを搬送して複数のノードが共有する資源を管理するためのコンピュータ読出可能媒体であって、前記命令の1つ以上のシーケンスが1つ以上のプロセッサによって実行されると、前記1つ以上のプロセッサは、

前記複数のノードのうち第1のノードの第1のキャッシュ内で前記資源を変更して前記資源の変更されたバージョンを作成するステップと、

前記第1のノードに障害が発生したときにどこで作業を開始するかを示す、前記第1のノードについてのチェックポイントを維持するステップと、

前記第1のキャッシュから永続的記憶装置に前記変更されたバージョンを最初に永続的に記憶することなしに、前記第1のキャッシュ内に前記変更されたバージョンの第1のコピーを保持する一方で、前記変更されたバージョンの第2のコピーを前記第1のキャッシュから前記複数のノードのうち第2のノードの第2のキャッシュに転送するステップと、

前記複数のノードのうち別のノードが、前記変更されたバージョンと少なくとも同じほど最近のものである前記資源のバージョンを永続的に記憶したことが示されることに応答

10

20

30

40

50

して、前記チェックポイントを進めるステップと実行する、コンピュータ読出可能媒体。

【請求項 4 2】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、
前記変更されたバージョンの前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで、前記第 1 のキャッシュ内に前記変更されたバージョンの少なくとも 1 つのコピーを保持するステップを実行するようにする命令をさらに含む、請求項 4 1 に記載のコンピュータ読出可能媒体。

【請求項 4 3】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、
前記複数のノードのうち別のノードが、前記変更されたバージョンと少なくとも同じほど最近のものである前記資源のバージョンを永続的に記憶したことが示されることに応答して、前記第 1 のノードで前記資源を開放するステップを実行するようにする命令をさらに含む、請求項 4 1 に記載のコンピュータ読出可能媒体。

10

【請求項 4 4】

命令の 1 つ以上のシーケンスを搬送して第 1 のキャッシュから第 2 のキャッシュへ資源を転送するためのコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサは、

前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから前記第 2 のキャッシュへと転送するステップと、

20

前記資源の第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 1 のコピーが前記第 1 のキャッシュ内で置換されることを防ぐステップとを実行する、コンピュータ読出可能媒体。

【請求項 4 5】

資源を転送するためのシステムであって、

第 1 のキャッシュを有するノードを含み、前記第 1 のキャッシュは、1 つ以上の他のノードに含まれる 1 つ以上の他のキャッシュのうち第 2 のキャッシュに通信可能に結合され、

前記ノードは、前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを第 1 のキャッシュから第 2 のキャッシュへと転送するように構成され、

30

前記ノードは、資源の前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 1 のキャッシュ内に資源の少なくとも 1 つのコピーを保持するように構成される、システム。

【請求項 4 6】

前記ノードは第 1 のデータベースサーバであり、前記 1 つ以上の他のノードのうち少なくとも 1 つのノードは前記第 2 のキャッシュを含む第 2 のデータベースサーバである、請求項 4 5 に記載のシステム。

【請求項 4 7】

40

前記ノードは、前記第 2 のコピーを前記第 2 のキャッシュに転送するよりも前に前記資源の前記第 1 のコピーが前記第 1 のキャッシュ内で変更されることを可能にするように構成され、

前記ノードは、前記第 2 のコピーを前記第 2 のキャッシュに転送した後に前記資源の前記第 1 のコピーが変更されることを防ぐように構成される、請求項 4 5 に記載のシステム。

【請求項 4 8】

前記ノードは、前記第 2 のコピーを前記第 2 のキャッシュに転送した後に、前記第 1 のコピーの開放の許可の要求をマスタノードに送信するように構成され、

前記ノードは、前記要求に応答して、前記マスタノードが前記第 1 のコピーまたはその

50

サクセサが永続的に記憶されるようにした後に、前記第 1 のコピーが開放可能であることを示すメッセージを前記マスタノードから受信するように構成される、請求項 4 5 に記載のシステム。

【請求項 4 9】

前記ノードは、前記第 1 のコピーの開放の許可の要求を送信するように構成される送信側プロセスを含み、

前記送信側プロセス以外のプロセスは前記資源の前記第 1 のコピーのサクセサを記憶する、請求項 4 8 に記載のシステム。

【請求項 5 0】

前記ノードは、前記第 1 のキャッシュ内に資源の少なくとも 1 つのコピーを、

前記第 1 のコピーを永続的に記憶しようとする前に、前記資源の永続的に記憶されたコピーが前記第 1 のコピーより最近のものであるかどうか決定し、

もし前記永続的に記憶されたコピーが前記第 1 のコピーより最近のものであれば、前記第 1 のコピーを永続的に記憶することなしに前記第 1 のコピーを開放し、

もし前記永続的に記憶されたコピーが前記第 1 のコピーより最近のもでなければ、前記第 1 のコピーを永続的に記憶することにより、保持するように構成される、請求項 4 5 に記載のシステム。

【請求項 5 1】

前記ノードは、変更許可を、第 1 のキャッシュに関連付けられる送信側プロセスから第 2 のキャッシュに関連付けられる受信側プロセスに、前記資源の前記第 2 のコピーとともに転送するように構成される、請求項 4 7 に記載のシステム。

【請求項 5 2】

前記資源にアクセスする許可はマスタノードによって管理され、

前記ノードは、前記変更許可の前記受信側プロセスへの転送についての確認を前記第 1 のノードが前記マスタノードから受信するより前に、前記変更許可を受信側プロセスに転送するように構成される、請求項 5 1 に記載のシステム。

【請求項 5 3】

前記ノードは、変更許可を、前記第 1 のキャッシュに関連付けられる送信側プロセスから前記第 2 のキャッシュに関連付けられる受信側プロセスに、前記資源の前記第 2 のコピーとともに転送するように構成され、

前記資源へのアクセス許可はマスタによって管理され、

前記変更許可の前記受信側プロセスへの転送は、前記変更許可の前記受信側プロセスへの転送の確認を前記マスタが受信するより前に実行される、請求項 4 5 に記載のシステム。

【請求項 5 4】

前記ノードは前記第 1 のキャッシュに関連付けられる送信側プロセスを含み、前記送信側プロセスは、前記第 2 のキャッシュに関連付けられる受信側プロセスから前記資源の要求を受信したマスタノードからメッセージを受信するように構成され、

前記送信側プロセスは、前記マスタノードからの前記メッセージに応答して、前記第 2 のコピーを前記受信側プロセスに転送する、請求項 4 5 に記載のシステム。

【請求項 5 5】

前記ノードは、前記第 2 のコピーが前記第 2 のキャッシュに転送された後にロックマネージャにロックを要求するように構成され、前記ロックは前記資源をディスクに書込む許可は与えるが前記資源を変更する許可は与えず、

前記ロックマネージャは、前記第 1 のコピーと少なくとも同じほど最近である前記資源のバージョンを有するプロセスを選択し前記ロックを前記選択されたプロセスに与えて、前記選択されたプロセスが前記資源の前記バージョンをディスクに書込むようにする、請求項 4 5 に記載のシステム。

【請求項 5 6】

前記資源の前記バージョンがディスクに書込まれるのに応答して、前記ロックマネージャ

10

20

30

40

50

ャが、前記バージョンよりも古い前記資源のすべてのバージョンが開放されるようにする、請求項 55 に記載のシステム。

【請求項 57】

前記資源のダーティコピーを保持するキャッシュの障害の後に、障害の発生したキャッシュが資源の最新バージョンを保持していたかどうかを決定するように構成されるマスタノードを含み、

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していたならば、マスタノードは、資源の最新パストイメージをディスクに書込み、資源の以前のすべてのパストイメージを開放し、前記障害の発生したキャッシュの復旧ログを適用して資源の最新バージョンを再構築するように構成される、請求項 45 に記載のシステム。

10

【請求項 58】

前記マスタノードは、もし前記障害の発生したキャッシュが資源の最新バージョンを保持していなかったならば、資源の最新バージョンをディスクに書込み資源のすべてのパストイメージを開放するように構成される、請求項 57 に記載のシステム。

【請求項 59】

前記資源のダーティバージョンを保持する複数のキャッシュの障害の後に、前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定し、もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、前記障害の発生したキャッシュの復旧ログをマージし適用して資源の最新バージョンを再構築するように構成されるマスタノードをさらに含む、請求項 45 に記載のシステム。

20

【請求項 60】

1 つ以上の共有されるディスクからの資源のダーティバージョンを保持する複数のキャッシュの間でデータを管理するためのシステムであって、

前記複数のキャッシュのうち 1 つのキャッシュを含むノードを含み、前記複数のキャッシュの各キャッシュに対して別々の復旧ログが維持され、前記複数のキャッシュのうち別のキャッシュが前記キャッシュ内に存在する資源のダーティバージョンを要求し、前記ノードは、前記ダーティバージョンを永続的に記憶することなしに前記ダーティバージョンを他のキャッシュに転送するように構成され、前記システムはさらに、

前記複数のキャッシュのうち障害の発生したキャッシュを、前記障害の発生したキャッシュに関連付けられた復旧ログに基づいて、前記複数のキャッシュのうち他のキャッシュの前記別々の復旧ログを検査することなしに、復旧するように構成されたマスタノードを含む、システム。

30

【請求項 61】

前記複数のキャッシュの各キャッシュは、複数のデータベースサーバの別々のデータベースサーバによって維持されるキャッシュである、請求項 60 に記載のシステム。

【請求項 62】

前記ダーティバージョンが存在するキャッシュは第 1 のキャッシュであり、

前記第 1 のキャッシュ内の資源のダーティバージョンは資源の第 1 のコピーであり、

前記ダーティバージョンを要求するキャッシュは第 2 のキャッシュであり、

ノードは、資源の第 2 のコピーを第 2 のキャッシュに転送することにより資源のダーティバージョンを転送するように構成され、

40

第 1 のノードは、前記第 2 のコピーを前記第 2 のキャッシュに転送するよりも前に前記資源の前記第 1 のコピーが前記第 1 のキャッシュ内で変更されることを可能にし、

前記第 1 のノードは、前記第 2 のコピーを前記第 2 のキャッシュに転送した後に前記資源の前記第 1 のコピーが変更されることを防ぐように構成される、請求項 60 に記載のシステム。

【請求項 63】

複数のノードによって共有される資源を管理するためのシステムであって、

第 1 のキャッシュを有するノードを含み、前記第 1 のキャッシュは、1 つ以上の他のノ

50

ードに含まれる 1 つ以上の他のキャッシュのうちの第 2 のキャッシュに通信可能に結合され、

前記ノードは、

前記ノードの第 1 のキャッシュ内で前記資源を変更して前記資源の変更されたバージョンを作成し、

前記ノードに障害が発生したときにどこで作業を開始するかを示す、前記ノードについてのチェックポイントを維持し、

前記第 1 のキャッシュから永続的記憶装置に前記変更されたバージョンを最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に前記変更されたバージョンの第 1 のコピーを保持する一方で、前記変更されたバージョンの第 2 のコピーを前記第 1 のキャッシュから第 2 のキャッシュに転送し、

10

前記複数のノードのうち別のノードが、前記変更されたバージョンと少なくとも同じほど最近のものである前記資源のバージョンを永続的に記憶したことが示されることに応答して、前記チェックポイントを進めるように構成される、システム。

【請求項 6 4】

前記ノードはさらに、

前記変更されたバージョンの前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで、前記第 1 のキャッシュ内に前記変更されたバージョンの少なくとも 1 つのコピーを保持するように構成される、請求項 6 3 に記載のシステム。

【請求項 6 5】

20

前記ノードはさらに、

前記複数のノードのうち別のノードが、前記変更されたバージョンと少なくとも同じほど最近のものである前記資源のバージョンを永続的に記憶したことが示されることに応答して、前記第 1 のノードで前記資源を開放するように構成される、請求項 6 3 に記載のシステム。

【請求項 6 6】

資源を転送するためのシステムであって、

第 1 のキャッシュを有するノードを含み、前記第 1 のキャッシュは、1 つ以上の他のノードに含まれる 1 つ以上の他のキャッシュのうち第 2 のキャッシュに通信可能に結合され、

30

前記ノードは、前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから前記第 2 のキャッシュへと転送するように構成され、

前記ノードは、資源の前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 1 のコピーが前記第 1 のキャッシュ内で置換されることを防ぐように構成される、システム。

【請求項 6 7】

資源のダーティバージョンを保持する 1 つ以上のキャッシュの障害の後に前記資源の復旧を管理するための方法であって、

40

障害の発生しなかった複数のキャッシュの各キャッシュにおいて資源のどのバージョンが保持されているのかを決定するステップと、

前記複数のキャッシュから、前記障害の発生しなかった複数のキャッシュ内の資源の他のいずれかのバージョンと少なくとも同じほど最近である資源の特定のバージョンを識別するステップと、

前記資源の特定のバージョンがディスクに書込まれるようにするステップとを含む、方法。

【請求項 6 8】

前記資源の特定のバージョンは資源の現在のコピーである、請求項 6 7 に記載の方法。

【請求項 6 9】

50

前記資源の特定のバージョンは資源のパストイメージコピーである、請求項 6 7 に記載の方法。

【請求項 7 0】

前記資源のパストイメージコピーは、資源の他のいずれかのパストイメージコピーと少なくとも同じほど最近のものである、請求項 6 9 に記載の方法。

【請求項 7 1】

前記資源のパストイメージコピーは、前記障害の発生しなかった複数のキャッシュに現在保持されている資源の他のいずれかのパストイメージコピーと少なくとも同じほど最近のものである、請求項 7 0 に記載の方法。

【請求項 7 2】

資源の最新バージョンがディスクに書込まれた後、前記複数のキャッシュが資源のパストイメージを開放するようにするステップをさらに含む、請求項 6 7 に記載の方法。

【請求項 7 3】

前記決定するステップ、識別するステップおよびディスクに書込まれるようにするステップは、ロックマネージャによって実行される、請求項 6 7 に記載の方法。

【請求項 7 4】

前記ロックマネージャは、前記 1 つ以上のキャッシュの障害の前に前記資源のマスタに指定された、生き残ったマスタである、請求項 7 3 に記載の方法。

【請求項 7 5】

前記ロックマネージャは、前記 1 つ以上のキャッシュの障害の後に前記資源のマスタに指定された、新しいマスタである、請求項 7 3 に記載の方法。

【請求項 7 6】

資源のダーティバージョンを保持する複数のキャッシュの障害の後に前記資源を復旧させるための方法であって、

前記障害の発生した複数のキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生した複数のキャッシュのいずれかが資源の最新バージョンを保持していたならば、

障害の発生したキャッシュのサブセットを決定するステップを含み、前記障害の発生したキャッシュのサブセットは、永続的に記憶された資源のバージョンに続いて前記資源を更新した、障害の発生したキャッシュのみを含み、さらに、

前記障害の発生したキャッシュのサブセットの復旧ログをマージし適用して資源の最新バージョンを再構築するステップを含む、方法。

【請求項 7 7】

命令の 1 つ以上のシーケンスを搬送して資源のダーティバージョンを保持する 1 つ以上のキャッシュの障害の後に前記資源の復旧を管理するためのコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサは、

障害の発生しなかった複数のキャッシュの各キャッシュにおいて資源のどのバージョンが保持されているのかを決定するステップと、

前記複数のキャッシュから、前記障害の発生しなかった複数のキャッシュ内の資源の他のいずれかのバージョンと少なくとも同じほど最近である資源の特定のバージョンを識別するステップと、

前記資源の特定のバージョンがディスクに書込まれるようにするステップとを実行する、コンピュータ読出可能媒体。

【請求項 7 8】

前記資源の特定のバージョンは資源の現在のコピーである、請求項 7 7 に記載のコンピュータ読出可能媒体。

【請求項 7 9】

前記資源の特定のバージョンは資源のパストイメージコピーである、請求項 7 7 に記載

10

20

30

40

50

のコンピュータ読出可能媒体。

【請求項 8 0】

前記資源のパストイメージコピーは、資源の他のいずれかのパストイメージコピーと少なくとも同じほど最近のものである、請求項 7 9 に記載のコンピュータ読出可能媒体。

【請求項 8 1】

前記資源のパストイメージコピーは、前記障害が発生しなかった複数のキャッシュに現在保持されている資源の他のいずれかのパストイメージコピーと少なくとも同じほど最近のものである、請求項 8 0 に記載のコンピュータ読出可能媒体。

【請求項 8 2】

1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、
資源の最新バージョンがディスクに書込まれた後、前記複数のキャッシュが資源のパストイメージを開放するようにするステップを実行するようにする命令をさらに含む、請求項 7 7 に記載のコンピュータ読出可能媒体。

10

【請求項 8 3】

前記命令の 1 つ以上のシーケンスが実行されると、前記 1 つ以上のプロセッサは、ロックマネージャに、前記決定するステップ、識別するステップおよび書込まれるようにするステップを実行させる、請求項 7 7 に記載のコンピュータ読出可能媒体。

【請求項 8 4】

前記ロックマネージャは、前記 1 つ以上のキャッシュの障害の前に前記資源のマスタに指定された、生き残ったマスタである、請求項 8 3 に記載のコンピュータ読出可能媒体。

20

【請求項 8 5】

前記ロックマネージャは、前記 1 つ以上のキャッシュの障害の後に前記資源のマスタに指定された、新しいマスタである、請求項 8 3 に記載のコンピュータ読出可能媒体。

【請求項 8 6】

命令の 1 つ以上のシーケンスを搬送して資源のダーティバージョンを保持する複数のキャッシュの障害の後に前記資源を復旧させるためコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサは、

前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定するステップを実行し、

30

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

障害の発生したキャッシュのサブセットを決定するステップを実行し、前記障害の発生したキャッシュのサブセットは、永続的に記憶された資源のバージョンに続いて前記資源を更新した、障害の発生したキャッシュのみを含み、さらに、

前記障害の発生したキャッシュのサブセットの復旧ログをマージし適用して資源の最新バージョンを再構築するステップを実行する、コンピュータ読出可能媒体。

【請求項 8 7】

資源のダーティバージョンを保持する 1 つ以上のキャッシュの障害の後に前記資源の復旧を管理するための装置であって、前記装置は、

40

障害の発生しなかった複数のキャッシュの各キャッシュにおいて資源のどのバージョンが保持されているのかを決定し、

前記複数のキャッシュから、前記障害の発生しなかった複数のキャッシュ内の資源の他のいずれかのバージョンと少なくとも同じほど最近である資源の特定のバージョンを識別し、

前記資源の特定のバージョンがディスクに書込まれるようにするよう構成される、装置。

【請求項 8 8】

前記資源の特定のバージョンは資源の現在のコピーである、請求項 8 7 に記載の装置。

【請求項 8 9】

50

前記資源の特定のバージョンは資源のパストイメージコピーである、請求項 8 7 に記載の装置。

【請求項 9 0】

前記資源のパストイメージコピーは、資源の他のいずれかのパストイメージコピーと少なくとも同じほど最近のものである、請求項 8 9 に記載の装置。

【請求項 9 1】

前記資源のパストイメージコピーは、前記障害の発生しなかった複数のキャッシュに現在保持されている資源の他のいずれかのパストイメージコピーと少なくとも同じほど最近のものである、請求項 9 0 に記載の装置。

【請求項 9 2】

前記装置はさらに、資源の最新バージョンがディスクに書込まれた後、前記複数のキャッシュが資源のパストイメージを開放するように構成される、請求項 8 7 に記載の装置。

【請求項 9 3】

前記装置はロックマネージャである、請求項 8 7 に記載の装置。

【請求項 9 4】

前記ロックマネージャは、前記 1 つ以上のキャッシュの障害の前に前記資源のマスタに指定された、生き残ったマスタである、請求項 9 3 に記載の装置。

【請求項 9 5】

前記ロックマネージャは、前記 1 つ以上のキャッシュの障害の後に前記資源のマスタに指定された、新しいマスタである、請求項 9 3 に記載の装置。

【請求項 9 6】

資源のダーティバージョンを保持する複数のキャッシュの障害の後に前記資源を復旧させるための装置であって、前記装置は、

前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定し、

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

障害の発生したキャッシュのサブセットを決定し、前記障害の発生したキャッシュのサブセットは、永続的に記憶された資源のバージョンに続いて前記資源を更新した、障害の発生したキャッシュのみを含み、さらに、

前記障害の発生したキャッシュのサブセットの復旧ログをマージし適用して資源の最新バージョンを再構築するように構成される、装置。

【請求項 9 7】

資源のダーティコピーを保持する第 1 のキャッシュの障害の後に前記資源を復旧させるための方法であって、

前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから第 2 のキャッシュへと転送するステップと、

前記障害の発生したキャッシュが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していたならば、前記障害の発生したキャッシュの復旧ログを資源の以前のバージョンに適用して資源の最新バージョンを再構築するステップとを含む、方法。

【請求項 9 8】

前記資源の以前のバージョンはディスク上に永続的に記憶される、請求項 9 7 に記載の方法。

【請求項 9 9】

前記資源の以前のバージョンは第 3 のキャッシュ内の資源のパストイメージである、請求項 9 7 に記載の方法。

【請求項 1 0 0】

前記第3のキャッシュ内の資源の前記パストイメーは、障害の発生しなかった複数のキャッシュのうちいずれかのキャッシュ内に現在保持されている資源の他のいずれかのパストイメーと少なくとも同じほど最新のものである、請求項99に記載の方法。

【請求項101】

前記障害の発生したキャッシュの復旧ログは、永続的記憶装置に前記資源の以前のバージョンを最初に永続的に記憶することなしに、前記第3のキャッシュ内の資源の以前のバージョンに適用される、請求項99に記載の方法。

【請求項102】

資源のダーティバージョンを保持する複数のキャッシュの障害の後に前記資源を復旧させるための方法であって、

10

第1のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第1のキャッシュ内に資源の第1のコピーを保持する一方で資源の第2のコピーを前記第1のキャッシュから第2のキャッシュへと転送するステップと、

前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

前記障害の発生したキャッシュの復旧ログをマージするステップと、

前記マージされた復旧ログを資源の以前のバージョンに適用して資源の最新バージョンを再構築するステップとを含む、方法。

20

【請求項103】

前記資源の以前のバージョンはディスク上に永続的に記憶される、請求項102に記載の方法。

【請求項104】

前記資源の以前のバージョンは、前記障害の発生した複数のキャッシュに含まれない第3のキャッシュ内の資源のパストイメーである、請求項102に記載の方法。

【請求項105】

前記第3のキャッシュ内の資源の前記パストイメーは、前記障害の発生した複数のキャッシュに含まれないいずれかのキャッシュ内に現在保持されている資源の他のいずれかのパストイメーと少なくとも同じほど最新のものである、請求項104に記載の方法。

30

【請求項106】

前記マージされた復旧ログは、永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第3のキャッシュ内の資源の以前のバージョンに適用される、請求項104に記載の方法。

【請求項107】

命令の1つ以上のシーケンスを搬送して資源のダーティコピーを保持する第1のキャッシュの障害の後に前記資源を復旧させるためのコンピュータ読出可能媒体であって、前記命令の1つ以上のシーケンスが1つ以上のプロセッサによって実行されると、前記1つ以上のプロセッサは、

前記第1のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第1のキャッシュ内に資源の第1のコピーを保持する一方で資源の第2のコピーを前記第1のキャッシュから第2のキャッシュへと転送するステップと、

40

前記障害の発生したキャッシュが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していたならば、前記障害の発生したキャッシュの復旧ログを資源の以前のバージョンに適用して資源の最新バージョンを再構築するステップとを実行する、コンピュータ読出可能媒体。

【請求項108】

前記資源の以前のバージョンはディスク上に永続的に記憶される、請求項107に記載のコンピュータ読出可能媒体。

50

【請求項 1 0 9】

前記資源の以前のバージョンは第 3 のキャッシュ内の資源のバストイメージである、請求項 1 0 7 に記載のコンピュータ読出可能媒体。

【請求項 1 1 0】

前記第 3 のキャッシュ内の資源の前記バストイメージは、障害の発生しなかった複数のキャッシュのうちいずれかのキャッシュ内に現在保持されている資源の他のいずれかのバストイメージと少なくとも同じほど最新のものである、請求項 1 0 9 に記載のコンピュータ読出可能媒体。

【請求項 1 1 1】

前記障害の発生したキャッシュの復旧ログは、永続的記憶装置に前記資源の以前のバージョンを最初に永続的に記憶することなしに、前記第 3 のキャッシュ内の資源の以前のバージョンに適用される、請求項 1 0 9 に記載のコンピュータ読出可能媒体。

10

【請求項 1 1 2】

命令の 1 つ以上のシーケンスを搬送し、資源のダーティバージョンを保持する複数のキャッシュの障害の後に前記資源を復旧させるためのコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサは、

第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから第 2 のキャッシュへと転送するステップと、

20

前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定するステップと、

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

前記障害の発生したキャッシュの復旧ログをマージするステップと、

前記マージされた復旧ログを資源の以前のバージョンに適用して資源の最新バージョンを再構築するステップとを実行する、コンピュータ読出可能媒体。

【請求項 1 1 3】

前記資源の以前のバージョンはディスク上に永続的に記憶される、請求項 1 1 2 に記載のコンピュータ読出可能媒体。

30

【請求項 1 1 4】

前記資源の以前のバージョンは、前記障害の発生した複数のキャッシュに含まれない第 3 のキャッシュ内の資源のバストイメージである、請求項 1 1 2 に記載のコンピュータ読出可能媒体。

【請求項 1 1 5】

前記第 3 のキャッシュ内の資源の前記バストイメージは、前記障害の発生した複数のキャッシュに含まれないいずれかのキャッシュ内に現在保持されている資源の他のいずれかのバストイメージと少なくとも同じほど最新のものである、請求項 1 1 4 に記載のコンピュータ読出可能媒体。

【請求項 1 1 6】

40

前記マージされた復旧ログは、永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 3 のキャッシュ内の資源の以前のバージョンに適用される、請求項 1 1 4 に記載のコンピュータ読出可能媒体。

【請求項 1 1 7】

資源のダーティコピーを保持する第 1 のキャッシュの障害の後に前記資源を復旧させるための装置であって、前記装置は、

前記第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから第 2 のキャッシュへと転送し、

前記障害の発生したキャッシュが資源の最新バージョンを保持していたかどうかを決定

50

し、

もし前記障害の発生したキャッシュが資源の最新バージョンを保持していたならば、前記障害の発生したキャッシュの復旧ログを資源の以前のバージョンに適用して資源の最新バージョンを再構築するように構成される、装置。

【請求項 1 1 8】

前記資源の以前のバージョンはディスク上に永続的に記憶される、請求項 1 1 7 に記載の装置。

【請求項 1 1 9】

前記資源の以前のバージョンは第 3 のキャッシュ内の資源のpastイメージである、請求項 1 1 7 に記載の装置。

10

【請求項 1 2 0】

前記第 2 のキャッシュ内の資源の前記pastイメージは、障害の発生しなかった複数のキャッシュのうちいずれかのキャッシュ内に現在保持されている資源の他のいずれかのpastイメージと少なくとも同じほど最新のものである、請求項 1 1 9 に記載の装置。

【請求項 1 2 1】

前記障害の発生したキャッシュの復旧ログは、永続的記憶装置に前記資源の以前のバージョンを最初に永続的に記憶することなしに、前記第 3 のキャッシュ内の資源の以前のバージョンに適用される、請求項 1 1 9 に記載の装置。

【請求項 1 2 2】

資源のダーティバージョンを保持する複数のキャッシュの障害の後に前記資源を復旧させるための装置であって、前記装置は、

20

第 1 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 1 のキャッシュ内に資源の第 1 のコピーを保持する一方で資源の第 2 のコピーを前記第 1 のキャッシュから第 2 のキャッシュへと転送し、

前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたかどうかを決定し、

もし前記障害の発生したキャッシュのいずれかが資源の最新バージョンを保持していたならば、

前記障害の発生したキャッシュの復旧ログをマージし、

前記マージされた復旧ログを資源の以前のバージョンに適用して資源の最新バージョンを再構築するように構成される、装置。

30

【請求項 1 2 3】

前記資源の以前のバージョンはディスク上に永続的に記憶される、請求項 1 2 2 に記載の装置。

【請求項 1 2 4】

前記資源の以前のバージョンは、前記障害の発生した複数のキャッシュに含まれない第 3 のキャッシュ内の資源のpastイメージである、請求項 1 2 2 に記載の装置。

【請求項 1 2 5】

前記第 3 のキャッシュ内の資源の前記pastイメージは、前記障害の発生した複数のキャッシュに含まれないいずれかのキャッシュ内に現在保持されている資源の他のいずれかのpastイメージと少なくとも同じほど最新のものである、請求項 1 2 4 に記載の装置。

40

【請求項 1 2 6】

前記マージされた復旧ログは、永続的記憶装置に前記資源を最初に永続的に記憶することなしに、前記第 3 のキャッシュ内の資源の以前のバージョンに適用される、請求項 1 2 4 に記載の装置。

【請求項 1 2 7】

複数のノードが使用する資源を管理するための方法であって、

前記複数のノードのうち第 1 のノードから資源の要求を受信するステップを含み、前記第 1 のノードは第 1 のキャッシュを含み、さらに、

前記複数のノードのうち第 2 のノードを識別するステップを含み、前記第 2 のノードは

50

資源の第 1 のコピーを有する第 2 のキャッシュを含み、さらに、

前記第 2 のノードが、前記第 2 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、資源の第 2 のコピーを前記第 2 のキャッシュから前記第 1 のキャッシュに転送するようにするステップと、

資源の前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 2 のキャッシュ内に資源の少なくとも 1 つのコピーが保持されるようにするステップとを含む、方法。

【請求項 1 2 8】

前記第 1 のノードは第 1 のデータベースサーバであり、前記第 2 のノードは第 2 のデータベースサーバである、請求項 1 2 7 に記載の方法。

10

【請求項 1 2 9】

前記複数のノードのうち第 3 のノードからさらに許可の要求を受信するステップをさらに含み、前記第 3 のノードは第 3 のキャッシュ内において資源の第 3 のコピーを含み、前記許可は前記第 3 のノードが前記資源をディスクに書込むことを許可するが前記第 3 のノードが前記資源を変更することは許可せず、さらに、

前記複数のノードのうち第 4 のノードを識別するステップを含み、前記第 4 のノードは第 4 のキャッシュ内において資源の第 4 のコピーを含み、前記第 4 のコピーは前記第 3 のコピーと少なくとも同じほど最近のものであり、さらに、

前記第 4 のノードが前記第 4 のコピーをディスクに書込むようにするステップを含む、請求項 1 2 7 に記載の方法。

20

【請求項 1 3 0】

前記第 4 のノードが前記第 4 のコピーをディスクに書込むことに応答して、前記第 4 のコピーよりも古い前記資源のすべてのコピーが開放されるようにするステップをさらに含む、請求項 1 2 9 に記載の方法。

【請求項 1 3 1】

前記第 4 のノードは前記第 3 のノードである、請求項 1 2 9 に記載の方法。

【請求項 1 3 2】

前記第 4 のノードは前記第 3 のノードでない、請求項 1 2 9 に記載の方法。

【請求項 1 3 3】

前記資源の前記第 4 のコピーは資源のパストイメージである、請求項 1 2 9 に記載の方法。

30

【請求項 1 3 4】

前記資源のパストイメージは、前記複数のノードが現在保持している資源の他のいずれかのパストイメージと少なくとも同じほど最近のものである、請求項 1 3 3 に記載の方法。

【請求項 1 3 5】

前記資源の前記第 4 のコピーは資源のカレントバージョンである、請求項 1 2 9 に記載の方法。

【請求項 1 3 6】

前記第 4 のノードが前記第 4 のコピーをディスクに書込むようにするステップは、

40

前記第 4 のノードに、前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可するさらなる許可を与えるステップを含み、前記さらなる許可は前記第 4 のノードが前記第 4 のコピーを変更することを許可しない、請求項 1 2 9 に記載の方法。

【請求項 1 3 7】

前記第 4 のノードに与えられる前記さらなる許可は、前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可する書込ロックであり、前記書込ロックは、前記第 4 のノードが前記第 4 のコピーを変更することは許可しない、請求項 1 3 6 に記載の方法。

【請求項 1 3 8】

前記第 4 のノードに与えられる前記さらなる許可は、前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可する書込トークンであり、前記書込トークンは、前記

50

第 4 のノードが前記第 4 のコピーを変更することは許可しない、請求項 1 3 6 に記載の方法。

【請求項 1 3 9】

前記第 4 のノードが前記第 4 のコピーをディスクに書込むようにするステップは、状態を前記第 4 のコピーに関連付けるステップを含み、前記状態は前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可する、請求項 1 2 9 に記載の方法。

【請求項 1 4 0】

前記要求を受信する前に、前記第 2 のノードに、前記第 2 のノードが前記第 1 のコピーを変更することを許可する第 1 の許可を与えるステップをさらに含み、前記第 1 の許可は前記第 2 のノードが前記第 1 のコピーをディスクに書込むことは許可せず、さらに、

10

前記第 2 のノードが前記第 2 のコピーを前記第 1 のキャッシュに転送するようにする前に、前記第 2 のノードに、前記第 2 のノードが前記第 2 のコピーを保持することを要求する第 2 の許可を与えるステップを含み、前記第 2 の許可は前記第 2 のノードが前記第 2 のコピーを変更することを許可せずかつ前記第 2 のノードが前記第 2 のコピーをディスクに書込むことを許可しない、請求項 1 2 7 に記載の方法。

【請求項 1 4 1】

前記第 2 のノードに与えられる前記第 1 の許可は、前記第 2 のノードが前記第 1 のコピーを変更することを許可する変更ロックであり、

前記変更ロックは前記第 2 のノードが前記第 1 のコピーをディスクに書込むことは許可せず、

20

前記第 2 の許可は前記第 2 のノードが前記第 2 のコピーを保持することを要求する保持ロックであり、

前記保持ロックは前記第 2 のノードが前記第 2 のコピーを変更することを許可せずかつ前記第 2 のノードが前記第 2 のコピーをディスクに書込むことを許可しない、請求項 1 4 0 に記載の方法。

【請求項 1 4 2】

前記第 2 のノードに与えられる前記第 1 の許可は、前記第 2 のノードが前記第 1 のコピーを変更することを許可する変更トークンであり、

前記変更トークンは前記第 2 のノードが前記第 1 のコピーをディスクに書込むことは許可せず、

30

前記第 2 の許可は前記第 2 のノードが前記第 2 のコピーを保持することを要求する保持トークンであり、

前記保持トークンは前記第 2 のノードが前記第 2 のコピーを変更することを許可せずかつ前記第 2 のノードが前記第 2 のコピーをディスクに書込むことを許可しない、請求項 1 4 0 に記載の方法。

【請求項 1 4 3】

前記要求を受信する前に、前記第 1 のコピーを、前記第 2 のノードが前記第 1 のコピーを変更することを許可しかつ前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可する第 1 の状態と関連付けるステップと、

前記第 2 のノードが前記第 2 のコピーを前記第 1 のキャッシュに転送するようにする前に、前記第 1 のコピーを、前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可する第 2 の状態と関連付けるステップとをさらに含み、前記第 2 の状態は前記第 2 のノードが前記第 1 のコピーを変更または上書きすることを許可しない、請求項 1 2 7 に記載の方法。

40

【請求項 1 4 4】

前記第 1 の状態は、前記第 2 のノードが前記第 1 のコピーを変更することを許可しかつ前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可する現在の状態であり、

前記第 2 の状態は、前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可するバストイメージ書込可能状態であり、前記バストイメージ書込可能状態は前記第 2

50

のノードが前記第 1 のコピーを変更または上書きすることを許可しない、請求項 1 4 3 に記載の方法。

【請求項 1 4 5】

前記第 2 のノードが変更許可を前記第 2 のノードから前記第 1 のノードに転送するようにするステップをさらに含む、請求項 1 2 7 に記載の方法。

【請求項 1 4 6】

前記第 2 のノードから、前記第 2 のノードが前記変更許可を前記第 1 のノードに転送したというメッセージを受信するステップをさらに含む、請求項 1 4 5 に記載の方法。

【請求項 1 4 7】

前記第 2 のノードが前記変更許可を前記第 1 のノードに転送したという前記メッセージは、前記第 2 のノードが前記変更許可を前記第 1 のノードに転送した後に受信される、請求項 1 4 6 に記載の方法。

10

【請求項 1 4 8】

前記第 1 のノードが前記変更許可を保持していることを示すデータを記憶するステップをさらに含む、請求項 1 4 6 に記載の方法。

【請求項 1 4 9】

複数のノードの共有資源を管理するための方法であって、

前記共有資源に変更許可を与えるのとは別に、前記共有資源に書込許可を与えるステップを含み、

前記書込許可を保持することは、前記書込許可を保持するノードが、ディスクに、前記書込許可を保持するノードのキャッシュに保持される共有資源のコピーを書込むことを許可し、前記書込許可を保持することは、前記書込許可を保持するノードが、前記書込許可を保持するノードのキャッシュに保持される共有資源のコピーを変更することは許可せず、

20

前記変更許可を保持することは、前記変更許可を保持するノードが、前記変更許可を保持するノードのキャッシュに保持される共有資源のコピーを変更することを許可し、前記変更許可を保持することは、前記変更許可を保持するノードが、前記変更許可を保持するノードのキャッシュに保持される共有資源のコピーを書込むことは許可しない、方法。

【請求項 1 5 0】

前記共有資源についての前記書込許可は前記複数のノードのうち 1 つのノードのみに時を選ばず与えられ、前記共有資源についての前記変更許可は前記複数のノードのうち 1 つのノードのみに時を選ばず与えられる、請求項 1 4 9 に記載の方法。

30

【請求項 1 5 1】

前記書込許可および前記変更許可は前記複数のノードのうち異なるノードによって保持される、請求項 1 4 9 に記載の方法。

【請求項 1 5 2】

前記書込許可および前記変更許可は前記複数のノードのうち同一のノードによって保持される、請求項 1 4 9 に記載の方法。

【請求項 1 5 3】

前記共有資源に前記書込許可を与えるのとは別にかつ前記共有資源に前記変更許可を与えるのとは別に、前記共有資源に保持許可を与えるステップをさらに含み、

40

前記保持許可を保持することは、前記保持許可を保持するノードが、前記保持許可を保持するノードのキャッシュに保持される資源のコピーを保持することを要求する、請求項 1 4 9 に記載の方法。

【請求項 1 5 4】

前記複数のノードのうち 1 組のノードの各ノードに対し、前記共有資源に前記書込許可を与えるのとは別にかつ前記共有資源に前記変更許可を与えるのとは別に、前記共有資源に保持許可を与えるステップをさらに含み、

前記保持許可を保持することは、前記保持許可を保持する 1 組のノードの各ノードが、前記保持許可を保持する 1 組のノードの各ノードのキャッシュに保持される共有資源のコ

50

ピーを保持することを要求する、請求項 1 4 9 に記載の方法。

【請求項 1 5 5】

前記保持許可を保持することは、前記保持許可を保持するノードが、前記保持許可を保持するノードのキャッシュに保持される共有資源のコピーまたはそのサクセサが永続的に記憶されるまで、前記保持許可を保持するノードのキャッシュに保持される共有資源のコピーを保持することを要求する、請求項 1 5 3 に記載の方法。

【請求項 1 5 6】

命令の 1 つ以上のシーケンスを搬送して複数のノードが使用する資源を管理するためのコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサは、

前記複数のノードのうち第 1 のノードから資源の要求を受信するステップを実行し、前記第 1 のノードは第 1 のキャッシュを含み、さらに、

前記複数のノードのうち第 2 のノードを識別するステップを実行し、前記第 2 のノードは資源の第 1 のコピーを有する第 2 のキャッシュを含み、さらに、

前記第 2 のノードが、前記第 2 のキャッシュから永続的記憶装置に前記資源を最初に永続的に記憶することなしに、資源の第 2 のコピーを前記第 2 のキャッシュから前記第 1 のキャッシュに転送するようにするステップと、

資源の前記第 1 のコピーまたはそのサクセサが永続的に記憶されるまで前記第 2 のキャッシュ内に資源の少なくとも 1 つのコピーが保持されるようにするステップとを実行する、コンピュータ読出可能媒体。

【請求項 1 5 7】

前記第 1 のノードは第 1 のデータベースサーバであり、前記第 2 のノードは第 2 のデータベースサーバである、請求項 1 5 6 に記載のコンピュータ読出可能媒体。

【請求項 1 5 8】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記複数のノードのうち第 3 のノードからさらに許可の要求を受信するステップを実行するようにする命令をさらに含み、前記第 3 のノードは第 3 のキャッシュ内において資源の第 3 のコピーを含み、前記許可は前記第 3 のノードが前記資源をディスクに書込むことを許可するが前記第 3 のノードが前記資源を変更することは許可せず、さらに、

前記複数のノードのうち第 4 のノードを識別するステップを実行するようにする命令を含み、前記第 4 のノードは第 4 のキャッシュ内において資源の第 4 のコピーを含み、前記第 4 のコピーは前記第 3 のコピーと少なくとも同じほど最近のものであり、さらに、

前記第 4 のノードが前記第 4 のコピーをディスクに書込むようにするステップを実行するようにする命令を含む、請求項 1 5 6 に記載のコンピュータ読出可能媒体。

【請求項 1 5 9】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記第 4 のノードが前記第 4 のコピーをディスクに書込むことに応答して、前記第 4 のコピーよりも古い前記資源のすべてのコピーが開放されるようにするステップを実行するようにする命令をさらに含む、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

【請求項 1 6 0】

前記第 4 のノードは前記第 3 のノードである、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

【請求項 1 6 1】

前記第 4 のノードは前記第 3 のノードでない、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

【請求項 1 6 2】

前記資源の前記第 4 のコピーは資源のパストイメージである、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

【請求項 1 6 3】

前記資源のパストイメージは、前記複数のノードが現在保持している資源の他のいずれ

10

20

30

40

50

かのパストイメージと少なくとも同じほど最近のものである、請求項 1 6 2 に記載のコンピュータ読出可能媒体。

【請求項 1 6 4】

前記資源の前記第 4 のコピーは資源のカレントバージョンである、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

【請求項 1 6 5】

前記第 4 のノードが前記第 4 のコピーをディスクに書込むようにするための命令は、前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記第 4 のノードに、前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可するさらなる許可を与えるステップを実行するようにする命令を含み、前記さらなる許可は前記第 4 のノードが前記第 4 のコピーを変更することを許可しない、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

10

【請求項 1 6 6】

前記第 4 のノードに与えられる前記さらなる許可は、前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可する書込ロックであり、前記書込ロックは、前記第 4 のノードが前記第 4 のコピーを変更することは許可しない、請求項 1 6 5 に記載のコンピュータ読出可能媒体。

【請求項 1 6 7】

前記第 4 のノードに与えられる前記さらなる許可は、前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可する書込トークンであり、前記書込トークンは、前記第 4 のノードが前記第 4 のコピーを変更することは許可しない、請求項 1 6 5 に記載のコンピュータ読出可能媒体。

20

【請求項 1 6 8】

前記第 4 のノードが前記第 4 のコピーをディスクに書込むようにするための命令は、前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

状態を前記第 4 のコピーに関連付けるステップを実行するようにする命令を含み、前記状態は前記第 4 のノードが前記第 4 のコピーをディスクに書込むことを許可する、請求項 1 5 8 に記載のコンピュータ読出可能媒体。

【請求項 1 6 9】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

30

前記要求を受信する前に、前記第 2 のノードに、前記第 2 のノードが前記第 1 のコピーを変更することを許可する第 1 の許可を与えるステップを実行するようにする命令をさらに含み、前記第 1 の許可は前記第 2 のノードが前記第 1 のコピーをディスクに書込むことは許可せず、さらに、

前記第 2 のノードが前記第 2 のコピーを前記第 1 のキャッシュに転送するようにする前に、前記第 2 のノードに、前記第 2 のノードが前記第 2 のコピーを保持することを要求する第 2 の許可を与えるステップを実行するようにする命令を含み、前記第 2 の許可は前記第 2 のノードが前記第 2 のコピーを変更することを許可せずかつ前記第 2 のノードが前記第 2 のコピーをディスクに書込むことを許可しない、請求項 1 5 6 に記載のコンピュータ読出可能媒体。

40

【請求項 1 7 0】

前記第 2 のノードに与えられる前記第 1 の許可は、前記第 2 のノードが前記第 1 のコピーを変更することを許可する変更ロックであり、

前記変更ロックは前記第 2 のノードが前記第 1 のコピーをディスクに書込むことは許可せず、

前記第 2 の許可は前記第 2 のノードが前記第 2 のコピーを保持することを要求する保持ロックであり、

前記保持ロックは前記第 2 のノードが前記第 2 のコピーを変更することを許可せずかつ前記第 2 のノードが前記第 2 のコピーをディスクに書込むことを許可しない、請求項 1 6 9 に記載のコンピュータ読出可能媒体。

50

【請求項 1 7 1】

前記第 2 のノードに与えられる前記第 1 の許可は、前記第 2 のノードが前記第 1 のコピーを変更することを許可する変更トークンであり、

前記変更トークンは前記第 2 のノードが前記第 1 のコピーをディスクに書込むことは許可せず、

前記第 2 の許可は前記第 2 のノードが前記第 2 のコピーを保持することを要求する保持トークンであり、

前記保持トークンは前記第 2 のノードが前記第 2 のコピーを変更することを許可せずかつ前記第 2 のノードが前記第 2 のコピーをディスクに書込むことを許可しない、請求項 1 6 9 に記載のコンピュータ読出可能媒体。

10

【請求項 1 7 2】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記要求を受信する前に、前記第 1 のコピーを、前記第 2 のノードが前記第 1 のコピーを変更することを許可しかつ前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可する第 1 の状態と関連付けるステップと、

前記第 2 のノードが前記第 2 のコピーを前記第 1 のキャッシュに転送するようにする前に、前記第 1 のコピーを、前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可する第 2 の状態と関連付けるステップとを実行するようにする命令をさらに含み、前記第 2 の状態は前記第 2 のノードが前記第 1 のコピーを変更または上書きすることを許可しない、請求項 1 5 6 に記載のコンピュータ読出可能媒体。

20

【請求項 1 7 3】

前記第 1 の状態は、前記第 2 のノードが前記第 1 のコピーを変更することを許可しかつ前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可する現在の状態であり、

前記第 2 の状態は、前記第 2 のノードが前記第 1 のコピーをディスクに書込むことを許可するバストイメージ書込可能状態であり、前記バストイメージ書込可能状態は前記第 2 のノードが前記第 1 のコピーを変更または上書きすることを許可しない、請求項 1 7 2 に記載のコンピュータ読出可能媒体。

【請求項 1 7 4】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記第 2 のノードが変更許可を前記第 2 のノードから前記第 1 のノードに転送するようにするステップを実行するようにする命令をさらに含む、請求項 1 5 6 に記載のコンピュータ読出可能媒体。

30

【請求項 1 7 5】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記第 2 のノードから、前記第 2 のノードが前記変更許可を前記第 1 のノードに転送したというメッセージを受信するステップを実行するようにする命令をさらに含む、請求項 1 7 4 に記載のコンピュータ読出可能媒体。

【請求項 1 7 6】

前記第 2 のノードが前記変更許可を前記第 1 のノードに転送したという前記メッセージは、前記第 2 のノードが前記変更許可を前記第 1 のノードに転送した後に受信される、請求項 1 7 5 に記載のコンピュータ読出可能媒体。

40

【請求項 1 7 7】

前記 1 つ以上のプロセッサによって実行されると、前記 1 つ以上のプロセッサが、

前記第 1 のノードが前記変更許可を保持していることを示すデータを記憶するステップを実行するようにする命令をさらに含む、請求項 1 7 5 に記載のコンピュータ読出可能媒体。

【請求項 1 7 8】

命令の 1 つ以上のシーケンスを搬送して複数のノードの共有資源を管理するためのコンピュータ読出可能媒体であって、前記命令の 1 つ以上のシーケンスが 1 つ以上のプロセッ

50

サによって実行されると、前記１つ以上のプロセッサは、

前記共有資源に変更許可を与えるのとは別に、前記共有資源に書込許可を与えるステップを実行し、

前記書込許可を保持することは、前記書込許可を保持するノードが、ディスクに、前記書込許可を保持するノードのキャッシュに保持される共有資源のコピーを書込むことを許可し、前記書込許可を保持することは、前記書込許可を保持するノードが、前記書込許可を保持するノードのキャッシュに保持される共有資源のコピーを変更することは許可せず、

前記変更許可を保持することは、前記変更許可を保持するノードが、前記変更許可を保持するノードのキャッシュに保持される共有資源のコピーを変更することを許可し、前記変更許可を保持することは、前記変更許可を保持するノードが、前記変更許可を保持するノードのキャッシュに保持される共有資源のコピーを書込むことは許可しない、コンピュータ読出可能媒体。

10

【請求項 179】

前記共有資源についての前記書込許可は前記複数のノードのうち１つのノードのみに時を選ばず与えられ、前記共有資源についての前記変更許可は前記複数のノードのうち１つのノードのみに時を選ばず与えられる、請求項 178 に記載のコンピュータ読出可能媒体。

【請求項 180】

前記書込許可および前記変更許可は前記複数のノードのうち異なるノードによって保持される、請求項 178 に記載のコンピュータ読出可能媒体。

20

【請求項 181】

前記書込許可および前記変更許可は前記複数のノードのうち同一のノードによって保持される、請求項 178 に記載のコンピュータ読出可能媒体。

【請求項 182】

前記１つ以上のプロセッサによって実行されると、前記１つ以上のプロセッサが、

前記共有資源に前記書込許可を与えるのとは別にかつ前記共有資源に前記変更許可を与えるのとは別に、前記共有資源に保持許可を与えるステップを実行するようにする命令をさらに含み、

前記保持許可を保持することは、前記保持許可を保持するノードが、前記保持許可を保持するノードのキャッシュに保持される資源のコピーを保持することを要求する、請求項 178 に記載のコンピュータ読出可能媒体。

30

【請求項 183】

前記１つ以上のプロセッサによって実行されると、前記１つ以上のプロセッサが、

前記複数のノードのうち１組のノードの各ノードに対し、前記共有資源に前記書込許可を与えるのとは別にかつ前記共有資源に前記変更許可を与えるのとは別に、前記共有資源に保持許可を与えるステップを実行するようにする命令をさらに含み、

前記保持許可を保持することは、前記保持許可を保持する１組のノードの各ノードが、前記保持許可を保持する１組のノードの各ノードのキャッシュに保持される共有資源のコピーを保持することを要求する、請求項 178 に記載のコンピュータ読出可能媒体。

40

【請求項 184】

前記保持許可を保持することは、前記保持許可を保持するノードが、前記保持許可を保持するノードのキャッシュに保持される共有資源のコピーまたはそのサクセサが永続的に記憶されるまで、前記保持許可を保持するノードのキャッシュに保持される共有資源のコピーを保持することを要求する、請求項 182 に記載のコンピュータ読出可能媒体。

【請求項 185】

複数のノードの共有資源を管理するためのシステムであって、

前記共有資源に変更許可を与えるのとは別に、前記共有資源に書込許可を与えるためのメカニズムを含み、

前記書込許可は、前記書込許可を保持するノードが、ディスクに、前記書込許可を保持

50

するノードのキャッシュに保持される資源のコピーを書込むことを許可し、前記書込許可は、前記書込許可を保持するノードが、前記書込許可を保持するノードのキャッシュに保持される資源のコピーを変更することは許可せず、

変更許可は、前記変更許可を保持するノードが、前記変更許可を保持するノードのキャッシュに保持される資源のコピーを変更することを許可し、前記変更許可は、前記変更許可を保持するノードが、前記変更許可を保持するノードのキャッシュに保持される資源のコピーを書込むことは許可しない、システム。

【請求項 186】

前記メカニズムはさらに、前記複数のノードに書込許可を 1 つのみ 1 度だけ与え、前記複数のノードに変更許可を 1 つのみ 1 度だけ与えるように構成される、請求項 185 に記載のシステム。

10

【請求項 187】

前記メカニズムはさらに保持許可を与えるように構成され、前記保持許可は、前記保持許可を保持するノードが、前記保持許可を保持するノードの第 3 のキャッシュに保持される資源の第 3 のコピーを保持することを要求する、請求項 185 に記載のシステム。

【請求項 188】

前記保持許可は、前記保持許可を保持するノードが、前記資源の第 3 のコピーまたはそのサクセサが永続的に記憶されるまで前記資源の第 3 のコピーを保持することを要求する、請求項 187 に記載のシステム。

【請求項 189】

20

前記メカニズムはさらに複数の保持許可を与えるように構成され、前記複数の保持許可の各保持許可は、前記複数のノードのうち 1 つのノードに与えられ、各保持許可は、前記複数の保持許可のうち 1 つを保持するノードが、前記複数の保持許可のうち 1 つを保持するノードのキャッシュに保持される資源の別のコピーを保持することを要求する、請求項 185 に記載のシステム。

【請求項 190】

前記複数の保持許可の各保持許可は、前記複数の保持許可のうち 1 つを保持するノードが、資源の他のコピーを、前記資源の他のコピーまたはそのサクセサが永続的に記憶されるまで保持することを要求する、請求項 189 に記載のシステム。

【発明の詳細な説明】

30

【0001】

【発明の分野】

この発明は、あるノードがデータストアからデータを要求するとき要求されたデータの最も最近のバージョンが別のノードのキャッシュ内にあるときことに関連付けられるペナルティを低減するための技術に関する。

【0002】

【発明の背景】

スケーラビリティを向上させるため、データベースシステムの中には、（各々が別個に稼動する）2 つ以上のデータベースサーバがディスクメディア上に記憶されるなど、共有の記憶装置に同時にアクセスすることを可能にするものがある。各データベースサーバは、ディスクブロックなどの、共有資源をキャッシュするためのキャッシュを有する。そのようなシステムをここではパラレルサーバシステムと呼ぶ。

40

【0003】

パラレルサーバシステムに関連付けられる問題の 1 つに、「ピング」と呼ばれるものの可能性がある。あるサーバのキャッシュ内にある資源のバージョンが異なったサーバのキャッシュに与えられなければならないとき、ピングは起きる。したがって、ピングが起きるのは、データベースサーバ A がそのキャッシュ内の資源 x を変更した後、データベースサーバ B が資源 x の変更を要求するときである。データベースサーバ A および B は、典型的には、異なったノード上で稼動するが、場合によっては同じノード上で稼動することもあり得る。

50

【 0 0 0 4 】

ピングを処理するアプローチの1つを、ここでは「ディスク介入」アプローチと呼ぶ。ディスク介入アプローチは、中間的記憶装置としてディスクを使用して2つのキャッシュ間で資源の最新バージョンを転送する。したがって、上記の例では、ディスク介入アプローチは、データベースサーバ1が資源Xのそのキャッシュバージョンをディスクに書込み、データベースサーバ2がこのバージョンをディスクからそのキャッシュへと検索することを必要とする。ディスク介入アプローチは、資源のサーバ間転送ごとに2ディスクI/Oを必要とするので、パラレルサーバシステムのスケーラビリティが制限される。具体的には、ピングを処理するために必要とされるディスクI/Oは、比較的不経済で時間がかかり、システムに加えられるデータベースサーバの数が多ければ多いほど、ピングの数も多

10

【 0 0 0 5 】

しかしながら、ディスク介入アプローチは、単一のデータベースサーバの障害からの比較的効率よい復旧を提供する、というのもそのような復旧が必要とするのは障害の発生したデータベースサーバの復旧（再実行）ログを適用するだけであるからである。障害の発生したデータベースサーバの再実行ログを適用すると、障害の発生したデータベースサーバ上のトランザクションが障害の発生したサーバのキャッシュ内の資源に加えた、かかわった変更はすべて確実に復旧される。復旧の間の再実行ログの使用は、1997年1月21日出願の「復旧可能オブジェクト内のキャッシングデータ」（“CACHING DATA IN RECOVERABLE OBJECTS”）と題する米国特許出願連続番号第08/784,611号に詳細に記

20

【 0 0 0 6 】

ディスク介入アプローチを採用するパラレルサーバシステムは、典型的には、資源アクセスおよび変更に関するグローバルな調停のすべてが分散ロックマネージャー（DLM）によって行なわれる場合のプロトコルを使用する。例示的DLMの動作は、1996年6月24日出願の「ロックキャッシングのための方法および装置」（“METHOD AND APPARATUS FOR LOCK CACHING”）と題する米国特許出願連続番号第08/669,689号に詳細に記載され、その内容はここに引用により援用される。

【 0 0 0 7 】

典型的な分散ロックマネージャースystemでは、任意の所与の資源に属する情報は、資源に対応するロックオブジェクト内に記憶される。各ロックオブジェクトは、単一のノードのメモリ内に記憶される。ロックオブジェクトが記憶されているノード上にあるロックマネージャーは、そのロックオブジェクトおよびそれがカバーする資源のマスタと呼ばれる。

30

【 0 0 0 8 】

ピングを処理するためにディスク介入アプローチを採用するシステムでは、ピングは、さまざまなロックが関係する通信においてDLMを必要とする。具体的には、データベースサーバ（「要求サーバ」）が資源のアクセスを必要とするとき、データベースサーバは、それが適切なモード、すなわち、読出の場合には共有され、書込の場合には排他的であるモードにロックされた所望の資源を有するかどうかをチェックする。もし要求データベースサーバが正しいモードにロックされた所望の資源を有していなければ、または、資源に全くロックがされていないならば、要求サーバは、資源のマスタに要求を送信して特定のモードのロックを獲得する。

40

【 0 0 0 9 】

要求データベースサーバによってなされた要求は、資源の現在の状態と競合することがある（たとえば、別のデータベースサーバが資源に対する排他的なロックを現在保持している可能性がある）。もし競合がなければ、資源のマスタは、ロックを許可し許可を登録する。競合の場合には、資源のマスタは、競合解決プロトコルを開始する。資源のマスタは、競合するロックを保持するデータベースサーバ（「ホルダ」）に、下位の互換性のあるモードにそのロックをダウングレードするよう命令する。

50

【0010】

不幸にも、もしホルダ（たとえば、データベースサーバA）が所望の資源の更新された（「ダーティ」）バージョンをそのキャッシュ内に現在持っていなければ、それはそのロックを即座にダウングレードできない。そのロックをダウングレードするために、データベースサーバAは、「ハードピング」プロトコルと呼ばれるものを経る。ハードピングプロトコルに従って、データベースサーバAは、ディスクに書込まれるべき更新に関連付けられる再実行ログを強制し、資源をディスクに書込み、そのロックをダウングレードし、データベースサーバAが完了したことをマスタに通知する。通知を受取ると、マスタは、ロック許可を登録し、要求されたロックが許可されたことを要求サーバに通知する。この時点で、要求サーバBは、ディスクからそのキャッシュ内に資源を読み出す。

10

【0011】

上述したとおり、ディスク介入アプローチによっては、あるデータベースサーバによって更新された資源（「ダーティ資源」）を別のデータベースサーバに直接発送することはできない。そのような直接発送は、復旧に関連する問題のために、実行可能性がないと考えられる。たとえば、資源がデータベースサーバAで変更されてから、データベースサーバBに直接発送されたと仮定する。データベースサーバBでも、資源は変更され、データベースサーバAに発送し返される。データベースサーバAで、資源は3度目に変更される。各サーバが、別のサーバに資源を送る前にすべての再実行ログをディスクに記憶することによって、受信側が先の変更可能となると仮定する。

【0012】

3度目の更新の後に、データベースサーバAがだめになったと仮定する。データベースサーバAのログは、穴のあいた資源への変更のレコードを含む。具体的には、サーバAのログは、データベースサーバBによってなされたこれらの変更を含んでいない。正確には、サーバBによってなされた変更は、データベースサーバBのログ内に記憶されている。この時点で、資源を復旧するために、2つのログは適用される前にマージされなければならない。このログマージ動作は、もし実現されれば、障害の発生しなかったデータベースサーバを含む、データベースサーバの総数に比例して時間および資源を必要とするであろう。

20

【0013】

上述したディスク介入アプローチは、障害の後の復旧ログのマージに関連付けられる問題を回避するが、簡単で効率のよい復旧を支持する定常状態の平行サーバシステムの性能にペナルティを科す。直接発送アプローチは、ディスク介入アプローチに関連付けられるオーバーヘッドを回避するが、障害の発生した場合に複雑で非スケーラブルな復旧動作を伴う。

30

【0014】

以上に基づいて、復旧動作の複雑性または持続時間を激しく増大させることなしに、ピングに関連付けられるオーバーヘッドを低減するためのシステムおよび方法を提供することが明らかに望まれる。

【0015】

【発明の概要】

最初に資源をディスクに書込むことなしに、あるデータベースサーバのキャッシュから別のデータベースサーバのキャッシュへと資源を転送するための方法および装置が提供される。データベースサーバ（要求者）が資源を変更したい場合、要求者は資源の現在のバージョンを要求する。現在のバージョンを有するデータベースサーバ（ホルダ）は現在のバージョンを要求者に直接発送する。バージョンを発送すると、ホルダは資源を変更する許可を失うが、メモリ内に資源のコピーを引続き保持する。資源の保持されたバージョンまたはその後のバージョンがディスクに書込まれると、ホルダは資源の保持されたバージョンを廃棄することができる。他の態様では、ホルダは保持されたバージョンを廃棄しない。サーバ障害の場合には、障害の発生したサーバの再実行ログ内の変更のあったすべての資源の先のコピーを、必要に応じて、障害の発生したサーバの再実行ログを適用するため

40

50

の開始点として使用する。この技術を用いて、単一のサーバ障害（障害の最もよくある形態）は、資源へアクセスしていなかったさまざまなデータベースサーバの復旧ログをマージする必要なしに、復旧される。

【0016】

この発明は、同じ参照番号が同様の要素を指している添付の図面に例として示されるが、これに限定されるものでない。

【0017】

【好ましい発明の詳細な説明】

ピングに関連付けられるオーバーヘッドを低減するための方法および装置を記載する。以下の記載では、説明のため、この発明を完全に理解するために、多くの具体的な詳細を述べる。しかしながら、この発明がこれらの具体的な詳細なしに実施可能であることは当業者には明らかであろう。他のデータベースサーバでは、この発明を不要にわかりにくくすることを避けるために、周知の構造および装置がブロック図の形で示される。

10

【0018】

機能概要

この発明のある局面に従うと、最初にディスクに記憶することなしに、データベースサーバ間で直接資源の更新されたバージョンを発送することによってピングを処理し、これによってディスク介入アプローチに関連付けられるI/Oオーバーヘッドを回避する。さらに、単一の場合の障害復旧に関連付けられる問題は、資源が別のキャッシュに転送されたとしても、変更された資源またはその何らかのサクセサがディスクに書込まれるまで資源の変更されたバージョンがキャッシュ内で置換されることを防ぐことによって、回避される。

20

【0019】

説明のため、キャッシュ内で置換不可能である資源のコピーを、ここでは「留められた」資源と呼ぶ。留められた資源を置換可能にする動作を、資源を「解放する」と呼ぶ。

【0020】

MおよびWロックアプローチ

この発明のある局面に従うと、資源に対する変更許可とディスクへの書込許可とは分離される。したがって、キャッシュからディスクへ資源の更新されたバージョンを書込む許可を有するデータベースサーバが、必ずしも資源を更新する許可を有するとは限らない。逆に、資源のキャッシュされたバージョンを変更する許可を有するデータベースサーバが、そのキャッシュされたバージョンをディスクに書込む許可を有するとは限らない。

30

【0021】

ある実施例に従うと、許可のこの分離は、特殊なロックを使用することによって実施される。具体的には、資源を変更する許可は、「M」ロックによって与えられるであろうし、ディスクに資源を書込む許可は、「W」ロックによって与えられるであろう。しかしながら、ここに記載するようなMロックおよびWロックの使用は、資源の転送されたバージョンが、その資源またはそのサクセサがディスクに書込まれるまでキャッシュ内で置換されることを防ぐためのメカニズムの1つにすぎないことが注目される。

【0022】

図2を参照すると、この発明のある実施例に従って、MロックおよびWロックを使用するデータシステムにおいてピングに回答して実行されるステップを示す。ステップ200で、資源を変更したいデータベースサーバは、資源のマスタ（すなわち、資源のロックを管理するデータベースサーバ）にMロックを要求する。ステップ202で、マスタは、資源のMロックを現在保持するデータベースサーバ（「ホルダ」）に、2つのサーバを接続する通信チャネル（「相互接続」）を介する直接転送によって、Mロックを資源のそのキャッシュされたバージョンとともに転送するよう命令する。

40

【0023】

ステップ204で、ホルダは、資源の現在のバージョンおよびMロックを要求者に送る。ステップ206で、ホルダは、Mロックの転送をマスタに通知する。ステップ208で、

50

マスタは、資源のロック情報を更新して要求者がMロックを現在保持していることを示す。

【0024】

P I 資源

Mロックのホルダは、必ずしもWロックを有しているとは限らず、このためそのキャッシュ内に含まれる資源のバージョンをディスクに書出す許可を有していない可能性がある。したがって、転送データベースサーバ（すなわち、Mロックを最後に保持していたデータベースサーバ）は、未来のある時点でそのバージョンをディスクに書出すよう要求される可能性があるので、資源のそのバージョンをダイナミックメモリに留め続ける。転送データベースサーバ内に留まる資源のバージョンは、もし受信データベースサーバが資源のそのコピーを変更すれば、古くなる。転送データベースサーバは、受信データベースサーバ（またはそのサクセサ）が資源をいつ変更するかわかっているとは限らないので、転送データベースサーバは資源のコピーを送信した時点から、その保持されたバージョンを「古い可能性のあるデータ」として扱う。資源のそのような古い可能性のあるバージョンを、ここではパストイメージ資源（P I 資源）と呼ぶ。

10

【0025】

P I 資源の解放

資源のキャッシュされたバージョンが解放された後、これは新しいデータで上書される可能性がある。典型的には、資源のダーティバージョンは、資源をディスクに書込むことによって解放され得る。しかしながら、キャッシュ内にP I 資源を有するデータベースサーバが必ずしも、P I 資源をディスクに記憶する権利を持っているとは限らない。これらの状況下でP I 資源を解放するためのある技術が、図3に示される。

20

【0026】

図3を参照すると、データベースサーバがそのキャッシュ内のP I 資源を解放したいとき、これはWロックの要求を分散ロックマネージャ（DLM）に送信する。ステップ302で、DLMは次に、要求データベースサーバ、または、資源のより新しいバージョン（サクセサ）をそのキャッシュ内に有する何らかのデータベースサーバに、資源をディスクに書出すよう命令する。こうして、資源をディスクに書込むよう命令されたデータベースサーバは、ダブルロックを許可される。ダブルロックを許可されたデータベースサーバが資源をディスクに書込んだ後、データベースサーバはWロックを解放する。

30

【0027】

次に、DLMは、すべてのデータベースサーバにメッセージを送信して書出された資源のバージョンを示し（ステップ304）、この結果、資源のこれ以前のP I バージョンはすべて解放可能となる（ステップ306）。たとえば、ディスクに書込まれたバージョンが時間T10で変更されたと仮定する。それより前の時間T5で最後に変更された資源のバージョンを有するデータベースサーバは、ここで、これが記憶されているバッファを他のデータのために使用することができるであろう。しかしながら、それより後の時間T11で変更されたバージョンを有するデータベースサーバは、資源のそのバージョンをそのメモリ内に保持し続けなければならないであろう。

【0028】

MおよびWロックアプローチの下でのピング管理

この発明のある実施例に従って、図1を参照して記載するように、MおよびWロックアプローチを実現化してピングを処理してもよい。図1を参照すると、4つのデータベースサーバA、B、CおよびDのブロック図が示され、これらのサーバはすべて特定の資源を含むデータベースへのアクセスを有している。例示される時点では、データベースサーバA、BおよびCはすべて、資源のバージョンを有する。データベースサーバAのキャッシュ内に保持されるバージョンは、資源の最も最近に変更されたバージョンである（時間T10で変更された）。データベースサーバBおよびCに保持されるバージョンは、資源のP I バージョンである。データベースサーバDは、資源のマスタである。

40

【0029】

50

この時点で、別のデータベースサーバ（「要求者」）が資源を変更したいと仮定する。要求者は、マスタに変更ロックを要求する。マスタは、要求者からの競合する要求のために、データベースサーバAにコマンドを送信してロック（「BAST」）をダウンコンバートする。ダウンコンバートコマンドに応答して、資源の現在のイメージ（クリーンまたはダーティのいずれでも）が、データベースサーバAから要求者に、資源を変更する許可とともに発送される。こうして発送された許可は、資源をディスクに書込む許可を含んでいない。

【0030】

データベースサーバAがMロックを要求者に送ると、データベースサーバAはそのMロックを「保持」ロック（「Hロック」）にダウングレードする。Hロックは、データベースサーバAが留められたPIコピーを保持していることを示す。Hロックの所有権は、オーナーにPIコピーをそのバッファキャッシュ内に維持することを強制するが、PIコピーをディスクに書込むいかなる権利もそのデータベースサーバに与えない。同じ資源に対して複数の同時的Hホルダがあり得るが、一度に資源の書込ができるデータベースサーバは1以下であり、したがって資源のWロックを保持することのできるデータベースサーバはたった1つである。

【0031】

資源を発送するより前に、データベースサーバAはログが確実に強制されるようにする（すなわち、データベースサーバAによって資源になされた変更について生成された復旧ログが永続的に記憶されるようにする）。変更許可を送ることによって、データベースサーバAは、資源を変更する自らの権利を失う。資源のコピー（発送の時点ではそうであったような）は、発送データベースサーバAになおも維持されている。資源の発送の後に、データベースサーバA内に保持される資源のコピーは、PI資源である。

【0032】

優遇書込

データベースサーバがダーティ資源を別のデータサーバに直接発送した後、資源の保持されたコピーは留められたPI資源となり、解放されるまでそのバッファを別の資源に使用することはできない。PI資源を含むバッファをここでは、PIバッファと呼ぶ。これらのバッファは、データベースサーバのキャッシュ内に有効な空間を占有しており、やがては他のデータのために再利用されなければならない。

【0033】

バッファキャッシュ内のPIバッファ（古くなったまたはチェックポイントされた）を置換するために、ここでは「優遇書込」と呼ぶ新しいディスク書込プロトコルを採用する。優遇書込プロトコルに従って、データベースサーバが資源をディスクに書込む必要があるとき、データベースサーバは要求をDLMに送信する。DLMは、ディスクに書込まれるべき資源のバージョンを選択し、選択されたバージョンを有するデータベースサーバを見つけ、書込要求を開始したデータベースサーバに代わって、そのデータベースサーバにディスクへの資源の書込をさせる。資源をディスクに実際に書込むデータベースサーバは、資源の最新の軌跡に依存して、書込を要求したデータベースサーバであっても、または、ほかの他のデータベースサーバであってもよい。

【0034】

資源の選択されたバージョンをディスクに書込むことによって、ディスクに書込まれた選択されたバージョンと同じ古さまたはそれよりも古い、クラスタのすべてのバッファキャッシュ内の資源のPIバージョンはすべて解放される。ディスクに書込まれるべきバージョンを選択するために使用される規準を、以下により詳細に記載する。しかしながら、選択されたバージョンは、マスタに知られている最新のPIバージョンか、または、資源のカレントバージョン（「CURR」）のいずれかであり得る。カレントバージョン以外のバージョンを選択する利点の1つは、この別のバージョンの選択によって現在のコピーが妨害されることなく変更可能となることである。

【0035】

P I 資源を保持しているデータベースサーバは、資源のWロックを獲得しているならば、そのP Iコピーを書出すことができる。資源の書込は、さまざまなデータベースサーバ間でのC U R R 資源イメージの移動から切離される。

【 0 0 3 6 】

効率的要因

資源を別のデータベースサーバに発送するたびにP Iコピーを書込む必要はない。したがって、資源を永続的に記憶する目的は、ディスクコピーを十分最近のものにしておくことと、バッファキャッシュ内の置換不可能な資源の数を妥当なものにしておくこととである。さまざまな要因が、上述した優遇書込プロトコルを採用するシステムの効率性を決定する。具体的には、

(1) ディスクにダーティ資源を書込むことによって起きるI / O動作を最低限にすることと、

(2) 資源のディスクバージョンを十分に現在のものにしておくことによって障害後の復旧動作を迅速化することと、

(3) 留められたP I 資源でバッファキャッシュがオーバーフローすることを防ぐことが望まれる。

【 0 0 3 7 】

第1の規準を最大化すると第2および第3の規準に否定的影響が及び、その逆もまたある。したがって、トレードオフが必要である。この発明のある実施例に従うと、総I O経費に対する制御と併せてチェックポイントのさまざまな技術(臨時的継続的チェックポイントと混合されたL R U) を組合せるセルフチューニングアルゴリズムを使用してもよい。

【 0 0 3 8 】

最新書込アプローチ

上述した優遇書込プロトコルの代替を、ここでは最新書込アプローチと呼ぶ。最新書込アプローチに従うと、すべてのデータベースサーバが、そのP I 資源をディスクに書込む許可を有する。しかしながら、そうする前に、データベースサーバは資源のディスクベースのコピーに対するロックを獲得する。ロックを獲得した後、データベースサーバは、ディスクバージョンを、これを書込みたいP I バージョンと比較する。もしディスクバージョンの方が古ければ、P I バージョンがディスクに書込まれる。もしディスクバージョンの方が新しければ、P I バージョンは廃棄されてもよく、それが占有していたバッファは再利用可能である。

【 0 0 3 9 】

優遇書込プロトコルと違って、最新書込アプローチは、データベースサーバが、自己のP I バージョンをディスクに書込むことによって、またはディスクバージョンの方がより新しいことを決定することによって、自己のP I バージョンを解放可能にする。しかしながら、最新書込アプローチは、ディスクベースのコピーのロックに対する競合を増大させ、優遇書込アプローチでは起きなかったであろうディスクI / O を招く可能性がある。

【 0 0 4 0 】

許可ストリング

典型的なD L M は、限られた数のロックモードを使用することによって資源へのアクセスを管理し、ここではモードは互換性があるか競合しているかのいずれかである。ある実施例に従うと、資源へのアクセスを管理するメカニズムは、ロックモードを異なった種の許可および義務の集合と代用するよう拡張される。許可および義務は、たとえば、資源を書込み、資源を変更し、キャッシュ内の資源を維持するなどの許可を含んでもよい。具体的な許可および義務を以下により詳細に記載する。

【 0 0 4 1 】

ある実施例に従うと、許可および義務は、許可ストリングに符号化される。多くの許可は資源自体ではなく資源のバージョンに相関するので、許可ストリングは資源バージョン数によって増大するであろう。もし2つの異なった許可ストリングが、資源の同じバージョン(たとえば、変更のための現在バージョンまたは書込のためのディスクアクセス)に

10

20

30

40

50

対する同じ排他的許可を要求するならば、これらは競合する。そうでなければこれらは互換性がある。

【0042】

許可転送を使用する同時実行性

上述したとおり、資源があるデータベースサーバで変更され、別のデータベースサーバによってさらなる変更を要求されると、マスタは、資源のカレントコピー（CURRコピー）を保持するデータベースサーバに、そのMロック（変更する権利）を資源のCURRコピーとともに他のデータベースサーバに送るように命令する。重要なことには、Mロックの要求はマスタに送信されるが、許可は何らかの他のデータベースサーバ（先のMロックホルダ）によってなされる。この三者間メッセージングモデルは、ロック要求が最初にアドレスされたロックマネージャーを含むデータベースサーバからロック要求に対する応答が期待される、従来の双方向通信とはかなり異なる。

10

【0043】

この発明のある実施例に従うと、資源のCURRコピーのホルダ（たとえば、データベースサーバA）がMロックを別のデータサーバに送ると、データベースサーバAは、Mロックが転送されたことをマスタに通知する。しかしながら、データベースサーバAは、マスタが通知を受取ったという確認を待つことなく、そのような確認を受取る前にCURRコピーおよびMロックを送信する。待たないことによって、マスタとデータベースサーバAとの間の往復通信は転送に遅延をもたらすことなく、これによってプロトコルレイテンシがかなり節約される。

20

【0044】

許可は許可の現在ホルダから許可の要求者に直接転送されるので、マスタが常に、ロック許可の正確な全体像を知っているとは限らない。むしろ、マスタは、任意の所与の時間でのロックの正確な位置についてではなく、Mロックの軌跡についてのみ、「これを最近保持した」データベースサーバについてのみ知っている。ある実施例に従うと、この「レージーな」通知方式は、Mロックに適用可能であるが、Wロック、Xロック、またはSロック（またはその対応物）には適用可能でない。ロック方式のさまざまな実施例を以下により詳細に記載する。

【0045】

障害復旧

30

この発明のコンテキストにおいては、サーバに関連付けられるキャッシュがアクセス不可能となった場合、データベースサーバに障害が発生したという。ここに記載する技術を用いるダーティ資源の直接のサーバ間発送を採用するデータシステムは、単一サーバの障害に回答して復旧ログをマージする必要性を回避する。ある実施例に従って、単一のサーバの障害は、図4に示すとおり処理される。図4を参照して、単一のデータベースサーバに障害が発生すると、復旧プロセスは、障害の発生したデータベースサーバのキャッシュ内に保持される各資源について、以下のステップを実行する。

【0046】

（ステップ400） 資源の最新バージョンを保持するデータベースサーバを決定し、
（ステップ402） ステップ400で決定されたデータベースサーバが障害の発生したデータベースサーバでなければ、（ステップ404）決定されたデータベースサーバは資源のそのキャッシュされたバージョンをディスクに書込み、（ステップ406）資源のPIバージョンはすべて解放される。このバージョンは、資源に加えられた、かかわった変更（障害の発生したデータベースサーバによってなされたものを含む）を有するため、いかなるデータベースサーバの復旧ログも適用される必要がない。

40

【0047】

もしステップ402で決定されたデータベースサーバが障害の発生したデータベースサーバであれば、（ステップ408）資源の最新PIバージョンを保持するデータベースサーバは、資源のそのキャッシュされたバージョンをディスクに書出し、（ステップ410）先のPIバージョンはすべて解放される。ディスクに書出されたバージョンは、障害の発

50

生したデータベースサーバ以外のすべてのデータサーバによって資源に加えられたかかわった変更を有する。障害の発生したデータベースサーバの復旧ログを適用して（ステップ412）障害の発生したデータベースサーバによって加えられたかかわった変更を復旧する。

【0048】

代替的に、資源の最新PIバージョンを、ディスク上ではなくキャッシュ内の現在のバージョンを復旧するための開始点として使用してもよい。具体的には、障害の発生したデータベースサーバの復旧ログから適切なレコードを、キャッシュ内にある最新PIバージョンに直接適用して、最新PIバージョンを保持するデータベースサーバのキャッシュ内のカレントバージョンを再構築してもよい。

10

【0049】

複数のデータベースサーバの障害

複数のサーバ障害の場合に、最新PIコピーもいかなるCURRコピーも生き残らなかったとき、資源になされた変更が障害の発生したデータベースサーバの複数のログにわたって広がっていることが起こり得る。この状況下では、障害の発生したデータベースサーバのログはマージされなければならない。しかしながら、すべてのデータベースサーバのログではなく、障害の発生したデータベースサーバのログのみがマージされなければならない。したがって、復旧のために必要とされる作業量は、構成全体のサイズではなく障害の程度に比例する。

【0050】

20

どの障害の発生したデータベースサーバが資源を更新したかを決定することが可能なシステムにおいては、資源を更新した障害の発生したデータベースサーバのログのみがマージされ適用される必要がある。同様に、どの障害の発生したデータベースサーバが、資源の永続的に記憶されたバージョンの後に資源を更新したかを決定することのできるシステムにおいては、資源の永続的に記憶されたバージョンの後に資源を更新した、障害の発生したデータベースサーバのログのみがマージされ適用される必要がある。

【0051】

例示的動作

説明のために、例示的な一連の資源転送を図1を参照して記載する。一連の転送の間、資源は複数のデータベースサーバでアクセスされる。具体的には、資源がクラスタノードに沿って発送され変更されると、データベースサーバの1つでのチェックポイントによってこの資源の物理的I/Oが起こる。

30

【0052】

再び図1を参照すると、4つのデータベースサーバ、A、B、CおよびDがある。データベースサーバDが資源のマスタである。まずデータベースサーバCが資源を変更する。データベースサーバCは資源バージョン8を有する。この時点で、データベースサーバCは、この資源に対するMロック（排他的変更権）も有する。

【0053】

この時点で、データベースサーバBが、データベースサーバCが現在保持している資源を変更したいと仮定する。データベースサーバBは、資源のMロックの要求（1）を送信する。データベースサーバDは、資源に関連付けられるモディファイアキュー上に要求を置き、

40

- （a） 変更許可（Mロック）をデータベースサーバBに送り、
- （b） 資源の現在イメージをデータベースサーバBに送信し、
- （c） データベースサーバCのMロックをHロックにダウングレードするよう、データベースサーバCに命令する（メッセージ2：BAST）。

【0054】

このダウングレード動作の後に、Cは、そのバッファキャッシュ内に資源のそのバージョン（PIコピー）を維持させられる。

【0055】

50

データベースサーバCは、要求された動作を実行し、新しい変更に対してログをさらに強制してもよい。加えて、データベースサーバCは、これが動作を実行したこと（AST）をマスタにレージーに通知する（3AckM）。この通知は、データベースサーバCがバージョン8を維持していることもマスタに知らせる。データベースサーバCは、マスタからの確認を待たない。従って、データベースサーバBは、マスタがそれを知る前に、Mロックを得ることが可能である。

【0056】

一方で、データベースサーバAもまた資源を変更することを決定したとする。データベースサーバAは、メッセージ（4）をデータベースサーバDに送信する。このメッセージは、データベースサーバCからデータベースサーバDへの非同期の通知の前に、到着し得る

10

【0057】

データベースサーバD（マスタ）は、（Bがこれを得て変更した後に）資源をデータベースサーバAに送るよう、データベースサーバB、すなわちこの資源の最新と知られているモディファイアにメッセージ（5）を送信する。なお、データベースサーバDは、資源がそこにあるのかまだなのかを知らない。しかし、データベースサーバDは、資源がやがてBに到着することは知っている。

【0058】

データベースサーバBが資源を得て意図された変更をした後（現在Bは資源のバージョン9を有している）、これは自己のロックをHにダウングレードし、データベースサーバAに、資源のカレントバージョン（「CURR資源」）をMロックとともに送信する（6）。データベースサーバBはまた、レージーな通知（6AckM）をマスタに送信する。

20

【0059】

この資源はデータベースサーバAで変更されつつあるが、データベースサーバCでのチェックポイントメカニズムが、資源をディスクに書込むことを決定したとする。上記の非同期の事象に関しては、3AckMおよび6AckMの両方が既にマスタに到着していると仮定する。チェックポイント動作に応答して実行された動作を図5を参照して示す。

【0060】

図5を参照すると、データベースサーバCは、書込権限を含まない、バージョン8に対するHロックを保持しているので、データベースサーバCは、メッセージ1をマスタ（D）に送信してそのバージョンについてのW（書込）ロックを要求する。この時点ではもう、マスタは、（確認が到着したと仮定して）資源がデータベースサーバAに発送されたことを知っている。データベースサーバDは、資源書込の命令とともに、（非請求の）WロックをデータベースサーバAに送信する（2BastW）。

30

【0061】

一般的な場合においては、この命令は、送信通知が到着している最新のデータベースサーバへ（または、最新であると知られている送信者から資源を受取ると考えられるデータベースサーバへ）送られる。データベースサーバAは、資源のそのバージョンを書込む（3）。データベースサーバによって書込まれた資源は、資源のバージョン10である。このときまでに、もしさらなる要求者が資源を要求していれば、資源のカレントコピーはどこか他にあるであろう。ディスクは、書込が完了したとき確認する（4AckW）。

40

【0062】

書込が完了すると、データベースサーバAは、データベースサーバDに、バージョン10が現在ディスク上にあるという情報を与える（5AckW）。データベースサーバAは、（これは最初には要求していなかった）そのWロックを自発的にダウングレードする。

【0063】

マスタ（D）はデータベースサーバCに行って、要求されたWロックを許可する代わりに、書込が完成したことをCに通知する（6）。マスタは、現在のディスクバージョン数を全てのPIコピーのホルダに知らせ、これによってCでのこれ以前のPIコピーはすべて解放可能となる。このシナリオでは、データベースサーバCは、10より古いPIコピー

50

を有していないので、これはデータベースサーバCのロックをNULLにダウンコンバートする。

【0064】

マスタはまた、確認メッセージをデータベースサーバBに送信してデータベースサーバBに10より以前のそのPIコピーを解放するよう命令する(7AckW(10))。

【0065】

分散ロックマネージャ

従来のDLM論理と対照的に、ここに記載する直接発送技術を実現するシステムでのマスタは、データベースサーバでのロック状態について不完全な情報を有することがある。ある実施例に従うと、資源のマスタは、以下の情報およびデータ構造を維持する。

10

【0066】

(1) (変更または共有アクセスのいずれかのための)CURRコピー要求者のキュー(キューの長さの上限は、クラスタ内のデータベースサーバの数である)。このキューをここでは、カレント要求キュー(CQ)と呼ぶ。

【0067】

(2) 資源が別のCURR要求者に送信されると、送信側はレージーに(これが確認を待たないという意味では非同期に)マスタに事象について通知する。マスタは、最新のいくつかの送信者を追跡し続ける。これがCQ上のポインタである。

【0068】

(3) ディスク上の最新資源バージョンのバージョン数。

20

(4) Wロック許可およびW要求キュー。

【0069】

ある実施例に従うと、W許可は同期する。すなわち、これはマスタによってのみ許可され、マスタは、この資源についてのクラスタ内の書込要求者が1以下であることを確実にする。マスタが次の許可を出すことができるのは、先の書込が完了しWロックが解放されたと通知された後のみである。もし2以上のモディファイアがあれば、Wロックは書込の持続時間の間与えられ、書込の後に自動的に解放される。もしモディファイアが1つだけであれば、モディファイアはW許可を維持可能である。

【0070】

(5) そのそれぞれの資源バージョン数を備えるHロックホルダのリスト。これは、バッファキャッシュ内のPIコピーについての情報を(おそらく不完全であるが)与える。

30

【0071】

ディスクウォームアップ

ここに記載する直接発送は、資源のバッファキャッシュイメージとディスクイメージとのライフサイクルを大きく引き離すので、復旧の際にこのギャップを埋める必要がある。ある実施例に従うと、DLM復旧とバッファキャッシュ復旧との間に、復旧の新しいステップが加えられる。この新しい復旧ステップをここでは「ディスクウォームアップ」と呼ぶ。

【0072】

通常のキャッシュ動作の間、資源のマスタは、(キャッシュ復旧に先行する)DLM復旧の際に、資源の位置とPIコピーおよびCURRコピーの利用可能性とについておおよそしか知らないが、資源のマスタは、生き残ったデータベースサーバのバッファキャッシュ内の最新PIおよびCURRコピーの利用可能性について完全な情報を収集する。資源のマスタが、(もし障害より前に資源が障害の発生したデータベースサーバ上にマスタされていれば)新しいマスタであっても生き残ったマスタであっても、これは当てはまる。

40

【0073】

情報を収集した後、マスタは、どのデータベースサーバが資源の最新コピーを所有しているかを知る。「ディスクウォームアップ」段では、マスタは、資源のこの最新コピー(もし利用可能であればCURR、および、もしCURRコピーが障害の発生したデータベースサーバとともに消失していれば最新PIコピー)のオーナーにWロックを発行する。次

50

に、マスタは、このデータベースサーバに、資源をディスクに書込むよう命令する。書込が完了すると、すべての他のデータベースサーバは、そのHロックをNULLロックに変換する（なぜなら書込まれたコピーが最新の利用可能なものであるからである）。これらのロックがコンバートされた後、キャッシュ復旧は通常通り続行可能である。

【0074】

ディスクウォームアップ段の間、いくつかの最適化が可能である。たとえば、もし最新イメージが復旧を実行するデータベースサーバのパッファキャッシュ内にあれば、資源は必ずしもディスクに書込まれる必要はない。

【0075】

ロックベース方式の代替

データベースサーバ間での資源のダーティコピーを直接発送するためのさまざまな技術を、特殊なタイプのロック（Mロック、WロックおよびHロック）使用するロッキング方式をコンテキストとして記載した。具体的には、これらの特殊ロックを使用して、（1）資源のカレントバージョンを有するサーバのみが資源を変更することと、（2）資源の同じバージョンまたはより新しいバージョンがディスクに書込まれるまで、すべてのサーバが資源のそのPIバージョンを維持することと、（3）資源のディスクベースのバージョンが資源のより古いバージョンによって重ね書きされないこととを確実にする。

【0076】

しかしながら、ロックベースのアクセス制御方式は、この発明が実施可能であるコンテキストの1つにすぎない。たとえば、任意のさまざまなアクセス制御方式を用いてこれらの同じ3つの規則を実施してもよい。したがって、この発明は、特定のタイプのアクセス制御方式に限定されるものではない。

【0077】

たとえば、ロックをベースとして資源へのアクセスを管理する代わりに、アクセスは、トークンによって管理されてもよく、この場合各トークンが特定のタイプの許可を表わす。特定の資源のためのトークンが、上述した3つの規則が確実に実施されるように、パラレルサーバ間で転送されてもよい。

【0078】

同様に、規則は、状態ベースの方式を用いて実施され得る。状態ベースの方式では、資源のバージョンは、事象に応答して状態を変化させ、バージョンの状態がそのバージョンに対して実行可能である動作のタイプを決定する。たとえば、データベースサーバは、その「現在の」状態での資源のカレントバージョンを受取る。現在の状態は、資源の変更および資源のディスクへの書込を可能とする。データベースサーバが資源のカレントバージョンを別のノードに転送すると、保持されているバージョンは「PI書込可能」状態になる。PI書込可能状態では、バージョンは、（1）変更不可能であり、（2）量ね書き不可能であるが、（3）ディスクへの書込は可能である。資源の任意のバージョンがディスクに書込まれると、ディスクに書込まれたバージョンと同じまたはそれよりも古い、PI書込可能状態にあるバージョンのすべてが、「PI解放」状態に置かれる。PI解放状態では、バージョンは量ね書き可能であるが、ディスクへの書込または変更は不可能である。

【0079】

ハードウェア概要

図6は、この発明の実施例が実現可能であるコンピュータシステム600を示すブロック図である。コンピュータシステム600は、情報を受け渡すためのバス602または他の通信メカニズムと、バス602に結合され情報を処理するためのプロセッサ604とを含む。コンピュータシステム600はまた、ランダムアクセスメモリ（RAM）または他のダイナミック記憶装置などの主メモリ606を含み、これはバス602に結合されプロセッサ604によって実行されるべき命令および情報を記憶する。主メモリ606はまた、プロセッサ604によって実行されるべき命令の実行の間、一時的変数または他の中間情報を記憶するために使用され得る。コンピュータシステム600は、リードオンリメモリ（ROM）608または他の静的記憶装置をさらに含み、これはバス602に結合さ

10

20

30

40

50

れ静的情報およびプロセッサ604のための命令を記憶する。磁気ディスクまたは光学ディスクなどの記憶装置610が設けられこれはバス602に結合され情報および命令を記憶する。

【0080】

コンピュータシステム600は、バス602を介して陰極線管(CRT)などのディスプレイ612に結合されもよく、これはコンピュータユーザに情報を表示する。英数字および他のキーを含む入力デバイス614は、バス602に結合されプロセッサ604に情報およびコマンド選択を与える。ユーザ入力デバイスの別のタイプは、マウス、トラックボールまたはカーソル方向キーなどのカーソルコントロール616であって、これは方向情報およびコマンド選択をプロセッサ604に与え、かつ、ディスプレイ612上のカーソルの動きを制御する。この入力デバイスは典型的には、第1の軸(たとえばx)および第2の軸(たとえばy)の、2軸での2自由度を有し、これによってデバイスは画面での位置を特定することが可能となる。

10

【0081】

この発明は、ピングに関連付けられるオーバーヘッドを低減するコンピュータシステム600の使用に関する。この発明のある実施例に従うと、ピングに関連付けられるオーバーヘッドは、プロセッサ604が主メモリ606に含まれる1つ以上の命令の1つ以上のシーケンスを実行することに対応して、コンピュータシステム600によって低減される。そのような命令は、記憶装置610などの、別のコンピュータ読出可能媒体から主メモリ606に読出されてもよい。主メモリ606内に含まれる命令のシーケンスを実行することによって、プロセッサ604はここに記載するプロセスステップを実行する。代替の実施例では、ハードワイヤ回路をソフトウェア命令の代わりにまたはこれと組合せて使用してこの発明を実現してもよい。したがって、この発明の実施例は、ハードウェア回路およびソフトウェアの特定の組合せに限定されない。

20

【0082】

ここに用いる「コンピュータ読出可能媒体」という言葉は、プロセッサ604に命令を与えて実行させることに関与する任意の媒体を指す。そのような媒体は、不揮発性媒体、揮発性媒体および伝送媒体を含むがこれに限られるものではない、多くの形態を取ってもよい。不揮発性媒体は、たとえば、記憶装置610などの、光学ディスクまたは磁気ディスクを含む。揮発性媒体は、主メモリ606などの、ダイナミックメモリを含む。伝送媒体は、バス602を含むワイヤを含む、同軸ケーブル、銅線および光ファイバを含む。伝送媒体はまた、電波および赤外データ通信の間生成されるものなど、音波または光波の形態を取ってもよい。

30

【0083】

コンピュータ読出可能媒体の通常の形態は、たとえば、フロッピー、フレキシブルディスク、ハードディスク、磁気テープ、またはその他の磁気媒体、CD-ROM、その他の光学媒体、パンチカード、紙テープ、孔のパターンを備えるその他の物理的媒体、RAM、PROMおよびEPROM、FLASH-EPROM、その他のメモリチップまたはカートリッジ、以下に記載する搬送波、またはコンピュータが読出可能なその他の媒体の形態を含む。

40

【0084】

コンピュータ読出可能媒体のさまざまな形態は、1つ以上の命令の1つ以上のシーケンスをプロセッサ604に搬送して実行することにかかわり得る。たとえば、命令は最初に、遠隔コンピュータの磁気ディスク上に担持されてもよい。遠隔コンピュータは、命令をそのダイナミックメモリにロードし、モデムを使用して電話線を介して命令を送信することができる。コンピュータシステム600にローカルなモデムは、電話線上のデータを受信し、赤外送信器を使用してデータを赤外信号に変換することができる。赤外検出器は赤外信号で搬送されるデータを受信可能であり、適切な回路がデータをバス602上に与えることができる。バス602は、データを主メモリ606に搬送し、プロセッサ604はそこから命令を検索し実行する。主メモリ606によって受取られた命令は、プロセッサ6

50

04によって実行される前またはその後に、記憶装置610上にオプションとして記憶されてもよい。

【0085】

コンピュータシステム600は、1つ以上の記憶装置（たとえばディスクドライブ655）がコンピュータシステム600と1つ以上の他のCPU（たとえばCPU651）の両方にアクセス可能である、共有ディスクシステムに属する。例示のシステムでは、ディスクドライブ655への共有アクセスは、システムエリアネットワーク653によって与えられる。しかしながら、さまざまなメカニズムを代替的に使用して共有アクセスを与えてもよい。

【0086】

コンピュータシステム600はまた、バス602に結合される通信インターフェイス618を含む。通信インターフェイス618は、双方向のデータ通信を与え、これはネットワークリンク620に結合し、ネットワークリンクはローカルネットワーク622に接続される。たとえば、通信インターフェイス618は、統合サービスデジタル網（ISDN）カードまたはモデムであってもよく対応するタイプの電話線にデータ通信接続を与える。別の例として、通信インターフェイス618は、互換性のあるLANにデータ通信接続を与えるローカルエリアネットワーク（LAN）カードであってもよい。ワイアレスリンクが実現されてもよい。いかなるそのような実現化例でも、通信インターフェイス618は、さまざまなタイプの情報を表わすデジタルデータストリームを搬送する電気信号、電磁波信号または光学信号を送信し受信する。

【0087】

ネットワークリンク620は、典型的には、1つ以上のネットワークを介して他のデータデバイスにデータ通信を与える。たとえば、ネットワークリンク620は、ローカルネットワーク622を介してホストコンピュータ624またはインターネットサービスプロバイダ（ISP）626によって動作するデータ装置に接続してもよい。ISP626は、現在通常「インターネット」628と呼ばれるワールドワイドパケットデータ通信ネットワークを介して、データ通信サービスを提供する。ローカルネットワーク622とインターネット628とはどちらも、デジタルデータストリームを搬送する電気信号、電磁波信号または光学信号を使用する。さまざまなネットワークを通る信号と、ネットワークリンク620上および通信インターネット618を通る信号とは、デジタルデータをコンピュータシステム600へかつそこから搬送するものであるが、情報を転送する搬送波の例示的形態である。

【0088】

コンピュータシステム600は、ネットワーク、ネットワークリンク620および通信インターフェイス618を介して、プログラムコードを含め、メッセージを送信しデータを受信することが可能である。インターネットの例では、サーバ630は、インターネット628、ISP626、ローカルネットワーク622および通信インターフェイス618を介して、アプリケーションプログラムのために要求されたコードを伝送可能である。

【0089】

受信されたコードは、受信されたときにプロセッサ604によって実行されてもよいし、かつ/または記憶装置610または他の不揮発性装置に記憶されて後に実行されてもよい。このようにして、コンピュータシステム600は、搬送波の形でアプリケーションコードを獲得可能である。

【0090】

複数のデータベースサーバが共通の永続性記憶装置へのアクセスを有するときに生じるピングを参照してピングを処理するための技術が記載されるが、この技術はこのコンテキストに限定されるのではない。具体的には、これらの技術は、あるキャッシュに関連付けられるプロセスが現在のバージョンが他のキャッシュ内に位置する資源を要求する可能性のあるいかなる環境に適用されてもよい。そのような環境は、たとえば、異なったノード上のテキストサーバが同じテキスト材料へのアクセスを有するような環境、異なったノード

10

20

30

40

50

上のメディアサーバが同じビデオデータへのアクセスを有するような環境、などを含む。

【0091】

ここに記載する技術を用いてピングを処理すれば、資源のデータベースサーバ間の転送は効率よくなるので、動作可能時間性能は、データベースサーバの数およびデータベースサーバあたりのユーザの増加にあわせて増大する。加えて、この技術によって、データベースサーバの数の増加にあわせて増大する単一のデータベースサーバの障害（障害の最もよくあるタイプ）からの効率的な復旧が得られる。

【0092】

重要なことには、ここに記載する技術は、ディスク介入によってではなく、IPCトランスポートを介して資源を送信することによってピングを処理する。したがって、ピングを
10 もたらす、資源についてのディスクI/Oは、かなり解消される。同期I/Oを伴うのは、これがログ強制のために必要とされる場合においてのみである。加えて、ディスクI/Oはチェックポイントおよびバッファキャッシュ置換のために生じるが、そのようなI/Oは、クラスタにわたるバッファ発送を減速させることはない。

【0093】

ここに記載する直接発送技術はまた、ピングによって生じるコンテキスト切換の数を低減させるようにもなる。具体的には、プロトコルの関与者（要求者およびホルダ）とマスタとの間の往復メッセージのシーケンスは、要求者、マスタ、ホルダ、要求者からなる通信
20 トライアングルによって代用される。

【図面の簡単な説明】

20

【図1】 資源の最新バージョンのキャッシュからキャッシュへの転送を例示するブロック図である。

【図2】 この発明の実施例に従ってディスク介入なしにあるキャッシュから別のキャッシュへ資源を伝送するステップを例示するフローチャートである。

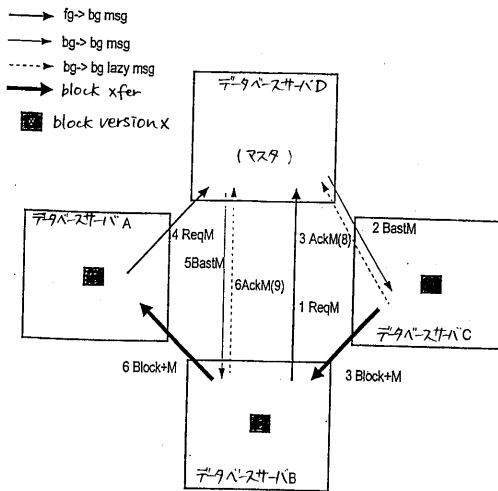
【図3】 この発明の実施例に従って、資源のラストイメージを解放するステップを例示するフローチャートである。

【図4】 この発明の実施例に従って単一のデータベースサーバの障害の後に復旧するステップを例示するフローチャートである。

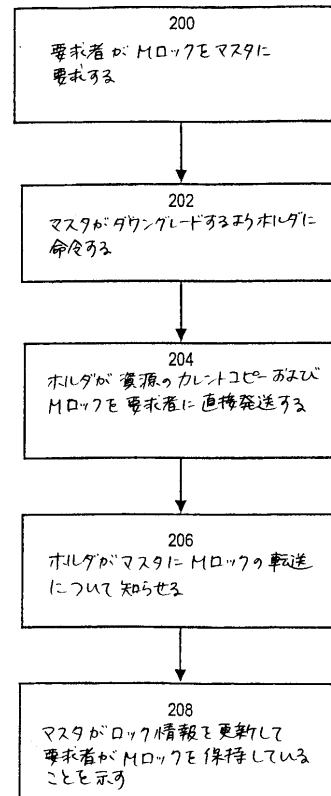
【図5】 この発明の実施例に従ってチェックポイントサイクルを例示するブロック図である。
30

【図6】 この発明の実施例が実現可能であるコンピュータシステムのブロック図である。
。

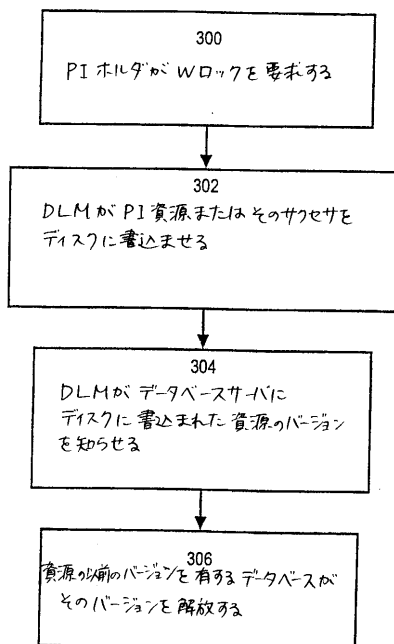
【図 1】



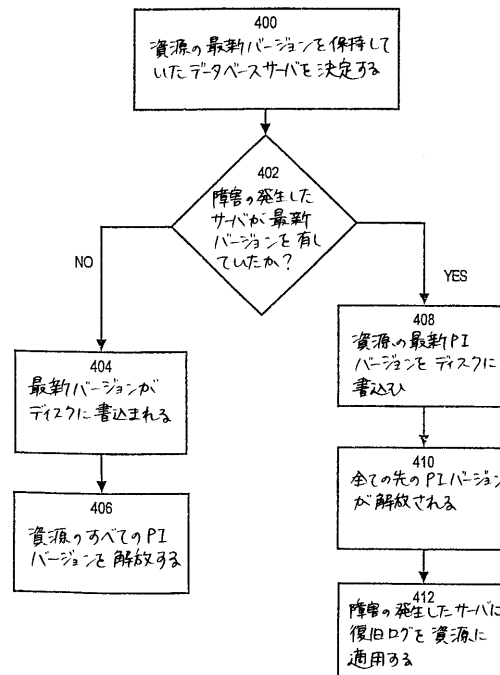
【図 2】



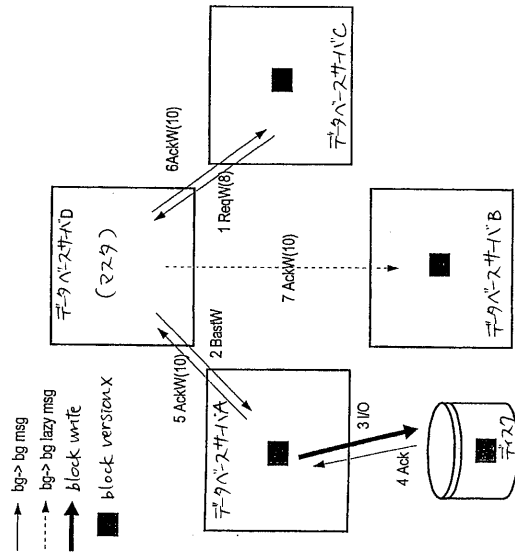
【図 3】



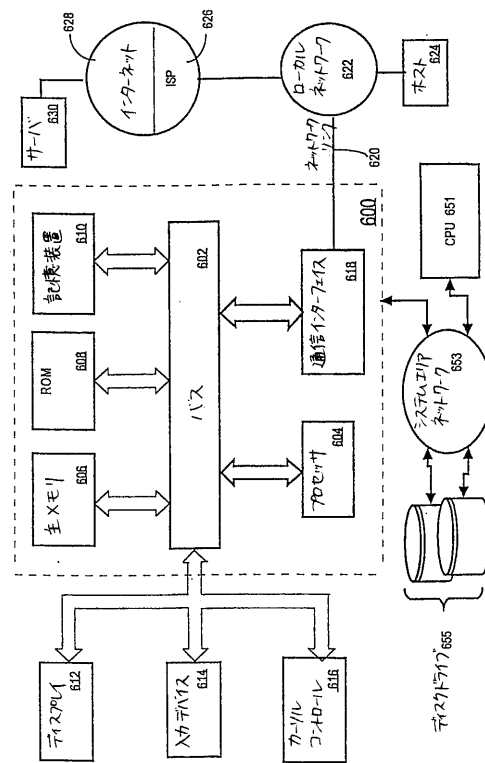
【図 4】



【図 5】



【図 6】



フロントページの続き

(74)代理人 100098316

弁理士 野田 久登

(74)代理人 100109162

弁理士 酒井 將行

(72)発明者 バンフォード, ロジャー・ジェイ

アメリカ合衆国、9 4 1 0 9 カリフォルニア州、サンフランシスコ、ハイド・ストリート、2 4
3 0

(72)発明者 クロッツ, ボリス

アメリカ合衆国、9 4 0 0 2 カリフォルニア州、ベルモント、ウィンディング・ウェイ、1 5 6
6

審査官 相崎 裕恒

(56)参考文献 特開平 0 4 - 8 4 3 5 0 (J P , A)

特開平 0 4 - 3 6 5 1 5 2 (J P , A)

特開平 1 1 - 8 5 6 0 0 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G06F 12/00,13/00

G06F 12/08-12

G06F 9/46/54

G06F 15/16-177