

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4472995号
(P4472995)

(45) 発行日 平成22年6月2日(2010.6.2)

(24) 登録日 平成22年3月12日(2010.3.12)

(51) Int. Cl. F I
G 0 6 F 12/00 (2006.01) G O 6 F 12/00 5 3 1 D
G 0 6 F 3/06 (2006.01) G O 6 F 3/06 3 0 4 F

請求項の数 18 (全 21 頁)

(21) 出願番号	特願2003-575277 (P2003-575277)	(73) 特許権者	500020944
(86) (22) 出願日	平成15年3月6日(2003.3.6)		マラソン テクノロジーズ コーポレイシ オン
(65) 公表番号	特表2005-519408 (P2005-519408A)		アメリカ合衆国 マサチューセッツ州 O 1 7 1 9 ボックスポロ マサチューセッ ツ アベニュー 1 3 0 0
(43) 公表日	平成17年6月30日(2005.6.30)	(74) 代理人	100082005
(86) 国際出願番号	PCT/US2003/006620		弁理士 熊倉 禎男
(87) 国際公開番号	W02003/077128	(74) 代理人	100067013
(87) 国際公開日	平成15年9月18日(2003.9.18)		弁理士 大塚 文昭
審査請求日	平成18年3月6日(2006.3.6)	(74) 代理人	100074228
(31) 優先権主張番号	10/090, 728		弁理士 今城 俊夫
(32) 優先日	平成14年3月6日(2002.3.6)	(74) 代理人	100086771
(33) 優先権主張国	米国 (US)		弁理士 西島 孝喜

最終頁に続く

(54) 【発明の名称】 増分不一致を使用するミラーコピーの生成方法

(57) 【特許請求の範囲】

【請求項 1】

コンピュータシステムにおいて、第1のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第1の記憶装置と関連し、第2のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第2の記憶装置と関連しており、該第1の記憶装置のミラーコピーを該第2の記憶装置に保持する方法であって、

前記第1のコントローラが前記第1の記憶装置において書き込み要求を受け取る段階と、

前記第1の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第1の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第1の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第1の記憶装置で受け取られた前記書き込み要求を前記第1のコントローラが処理する段階と、

10

前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第2の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第2の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第2の記憶装置において書き込み要求を前記第2のコントローラが受け取る段階と、

前記第1の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第1のコントローラが前記第2の記憶装置に周期的に送る段階と、

20

前記第2のコントローラが、前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、前記第2のコントローラが揮発性ディスク記憶装置において保存されている書き込み要求を前記第2の記憶装置の不揮発性ディスク記憶装置にコミットする段階と、

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第2の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込み要求に先行していたすべての書き込み要求に関連する不揮発性記憶装置書き込み対象データが前記第2の記憶装置の不揮発性記憶装置に書き込まれたことと、前記指定された書き込み要求に関連する不揮発性記憶装置書き込み対象データが前記第2の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第2の記憶装置によって前記第1のコントローラに送る段階と、を含む方法。

10

【請求項2】

前記同期コミットメッセージを受け取った後で、前記第2の記憶装置が前記指定された書き込み要求を処理する段階を更に含む請求項1に記載の方法。

【請求項3】

コンピュータシステムにおいて、第1のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第1の記憶装置と関連し、第2のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第2の記憶装置と関連しており、該第1の記憶装置のミラーコピーを該第2の記憶装置に保持する方法であって、

前記第1のコントローラが前記第1の記憶装置において書き込み要求を受け取る段階と

20

前記第1の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第1の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第1の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第1の記憶装置で受け取られた前記書き込み要求を前記第1のコントローラが処理する段階と、

前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第2の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第2の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第2の記憶装置において書き込み要求を前記第2のコントローラが受け取る段階と、

30

前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第2の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第2の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第2の記憶装置において受け取られた前記書き込み要求を前記第2のコントローラが処理する段階と、

前記第1の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第1のコントローラが前記第2の記憶装置に周期的に送る段階と、

前記第2のコントローラが、前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、前記第2のコントローラが揮発性ディスク記憶装置において保存されている書き込み要求を前記第2の記憶装置の不揮発性ディスク記憶装置にコミットする段階と、

40

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第2の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記第2の記憶装置の揮発性記憶装置の正常なキャッシュフラッシュを示す情報を含む確認メッセージを前記第2の記憶装置によって前記第1のコントローラに送る段階と、を含む方法。

【請求項4】

コンピュータシステムにおいて、第1のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第1の記憶装置と関連し、第2のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第2の記憶装置と関連しており、該第1の記憶装置のミラーコ

50

ピーを該第 2 の記憶装置に保持する方法であって、

前記第 1 のコントローラが前記第 1 の記憶装置において書き込み要求を受け取る段階と

、
前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 1 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 1 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 1 の記憶装置で受け取られた前記書き込み要求を前記第 1 のコントローラが処理する段階と、

前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 2 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 2 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 2 の記憶装置において書き込み要求を前記第 2 のコントローラが受け取る段階と、

前記第 1 の記憶装置によって処理された書き込み要求を識別する基準ラベルであって、書き込み要求に対して順次に割り当てられる基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第 1 のコントローラが前記第 2 の記憶装置に周期的に送る段階と、

前記第 2 のコントローラが、前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、前記第 2 のコントローラが揮発性ディスク記憶装置において保存されている書き込み要求を前記第 2 の記憶装置の不揮発性ディスク記憶装置にコミットする段階と、

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第 2 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込み要求に関連するデータが前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第 2 の記憶装置によって前記第 1 のコントローラに送る段階と、を含む方法。

【請求項 5】

前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に書き込み要求が保存される前に、前記同期コミットメッセージの基準ラベルによって識別された該書き込み要求に先立って発行された全ての書き込み要求を、前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に保存するように、前記第 2 のコントローラが書き込み要求を前記第 2 の記憶装置において受け取る段階は、前記第 2 のコントローラが要求をこの要求の基準ラベルによって順次に書き込む段階を含むことを特徴とする請求項 4 に記載の方法。

【請求項 6】

各記憶装置が基準ラベルと同じ順番で書き込み要求を受け取ることを特徴とする請求項 4 に記載の方法。

【請求項 7】

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に保存された書き込み要求によって影響を受ける記憶装置の領域を前記第 1 のコントローラが識別する段階を更に含む請求項 4 に記載の方法。

【請求項 8】

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置において保存されていた書き込み要求によって影響を受ける記憶装置の領域を前記第 1 のコントローラが識別する段階は、第 1 のビットマップ内に前記識別された記憶装置の領域を前記第 1 のコントローラが蓄積する段階を更に含む請求項 7 に記載の方法。

【請求項 9】

コンピュータシステムにおいて、第 1 のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第 1 の記憶装置と関連し、第 2 のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第 2 の記憶装置と関連しており、該第 1 の記憶装置のミラーコピーを該第 2 の記憶装置に保持する方法であって、

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存

10

20

30

40

50

するとともに、前記第 1 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 1 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 1 のコントローラが前記第 1 の記憶装置において書き込み要求を受け取る段階と、

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 1 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 1 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 1 の記憶装置で受け取られた前記書き込み要求を前記第 1 のコントローラが処理する段階と、

前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 2 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 2 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 2 の記憶装置において書き込み要求を前記第 2 のコントローラが受け取る段階と、

前記第 1 の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第 1 のコントローラが前記第 2 の記憶装置に周期的に送る段階と、

前記第 2 のコントローラが、前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、前記第 2 のコントローラが揮発性ディスク記憶装置において保存されている書き込み要求を前記第 2 の記憶装置の不揮発性ディスク記憶装置にコミットする段階と、

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第 2 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込み要求に関連する不揮発性記憶装置書き込み対象データが前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第 2 の記憶装置によって前記第 1 のコントローラに送る段階と、

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に保存された書き込み要求によって影響を受ける記憶装置の領域を前記第 1 のコントローラが識別する段階と、

前記第 1 のコントローラが前記識別された領域を第 1 のビットマップ中に蓄積する段階と、

前記同期コミットメッセージを送った後で、第 2 のビットマップ内に新規に識別された記憶装置の領域を前記第 1 のコントローラが蓄積する段階と、

前記処理された書き込み要求のデータが前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことを前記第 2 の記憶装置が確認した後で、前記書き込みデータが前記不揮発性記憶装置に正常に書き込まれたかどうかを示すステータスメッセージを前記第 2 のコントローラによって前記第 1 の記憶装置に送る段階と、

前記書き込みデータが正常に書き込まれたことを示す前記ステータスメッセージを受け取った後で、前記第 1 のビットマップを前記第 1 のコントローラが削除する段階と、を含む方法。

【請求項 10】

前記書き込みデータが正常に書き込まれなかったことを示すステータスメッセージを前記第 1 のコントローラが受け取った後で、前記第 2 のビットマップの内容を前記第 1 のビットマップに前記第 1 のコントローラがコピーし、前記第 2 のビットマップを前記第 1 のコントローラが削除する段階を更に含む請求項 9 に記載の方法。

【請求項 11】

前記第 1 のビットマップを削除した後で、前記第 2 のビットマップを前記第 1 のビットマップとして前記第 1 のコントローラが指定する段階を更に含む請求項 9 に記載の方法。

【請求項 12】

コンピュータシステムにおいて、第 1 のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第 1 の記憶装置と関連し、第 2 のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第 2 の記憶装置と関連しており、該第 1 の記憶装置のミラーコピーを該第 2 の記憶装置に保持する方法であって、

10

20

30

40

50

前記第1のコントローラが前記第1の記憶装置において書き込み要求を受け取る段階と

、
前記第1の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第1の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第1の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第1の記憶装置で受け取られた前記書き込み要求を前記第1のコントローラが処理する段階と、

前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第2の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第2の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記

10

第2の記憶装置において書き込み要求を前記第2のコントローラが受け取る段階と、
前記第1の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第1のコントローラが前記第2の記憶装置に周期的に送る段階と、

前記第2のコントローラが、前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、前記第2のコントローラが揮発性ディスク記憶装置において保存されている書き込み要求を前記第2の記憶装置の不揮発性ディスク記憶装置にコミットする段階と、

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第2の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込み要求に関連する不揮発性記憶装置書き込み対象データが前記第2の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第2の記憶装置によって前記第1のコントローラに送る段階と、

20

前記第2の記憶装置が書き込み要求を処理することができなかつた期間の後で、

前記第1のコントローラが前記第1のビットマップの内容を回復ビットマップにコピーする段階と、

前記回復ビットマップを使用して、前記第1の記憶装置から前記第2の記憶装置にコピーされることになる前記第1の記憶装置の記憶領域を前記第1のコントローラが識別する段階と、

前記第1の記憶装置の識別された記憶領域を前記第1のコントローラが前記第2の記憶装置にコピーする段階と、

30

第3のビットマップ内に前記第1の記憶装置において新規に受け取られた書き込み要求を前記第1のコントローラが蓄積する段階と、
を含む方法。

【請求項13】

コンピュータシステムにおいて、第1のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第1の記憶装置と関連し、第2のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第2の記憶装置と関連しており、該第1の記憶装置のミラーコピーを該第2の記憶装置に保持する方法であって、

前記第1のコントローラが前記第1の記憶装置において書き込み要求を受け取る段階と

40

、
前記第1の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第1の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第1の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第1の記憶装置で受け取られた前記書き込み要求を前記第1のコントローラが処理する段階と、

前記第2の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第2の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第2の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第2の記憶装置において書き込み要求を前記第2のコントローラが受け取る段階と、

50

前記第 1 の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第 1 のコントローラが前記第 2 の記憶装置に周期的に送る段階と、

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第 2 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込み要求に関連する不揮発性記憶装置書き込み対象データが前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第 2 の記憶装置によって前記第 1 のコントローラに送る段階と、

前記第 2 の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である第 2 の書き込み要求を指定する情報と共に第 2 の同期コミットメッセージを前記第 2 のコントローラが前記第 1 の記憶装置に送る段階と、

前記第 2 の同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第 1 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された第 2 の書き込み要求に関連する不揮発性記憶装置書き込み対象データが前記第 1 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第 1 の記憶装置によって前記第 2 のコントローラに送る段階と、
を含む方法。

【請求項 14】

前記処理された書き込み要求が前記第 1 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された第 2 の書き込み要求に関連するデータが前記第 1 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第 1 の記憶装置が前記第 2 のコントローラに送る段階は、前記指定された第 2 の書き込み要求に先行する全ての書き込み要求に関連するデータが前記第 1 の記憶装置の不揮発性記憶装置に書き込まれているデータである確認メッセージを前記第 1 の記憶装置が前記第 2 のコントローラに送る段階を含む請求項 13 に記載の方法。

【請求項 15】

コンピュータシステムにおいて、第 1 のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第 1 の記憶装置と関連し、第 2 のコントローラが揮発性記憶装置及び不揮発性記憶装置を具備する第 2 の記憶装置と関連しており、該第 1 の記憶装置のミラーコピーを該第 2 の記憶装置に保持する方法であって、

前記第 1 のコントローラが前記第 1 の記憶装置において書き込み要求を受け取る段階と、

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 1 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 1 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 1 の記憶装置で受け取られた前記書き込み要求を前記第 1 のコントローラが処理する段階と、

前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 2 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 2 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 2 の記憶装置において書き込み要求を前記第 2 のコントローラが受け取る段階と、

前記第 1 の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である書き込み要求を指定する情報と共に同期コミットメッセージを前記第 1 のコントローラが前記第 2 の記憶装置に周期的に送る段階と、

前記第 2 のコントローラが、前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、前記第 2 のコントローラが揮発性ディスク記憶装置において保存されている書き込み要求を前記第 2 の記憶装置の不揮発性ディスク記憶装置にコミットする段階と、

前記同期コミットメッセージを受け取った後で、前記処理された書き込み要求が前記第 2 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込

10

20

30

40

50

み要求に関連するデータが前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを前記第 2 の記憶装置によって前記第 1 のコントローラに送る段階と、

前記同期コミットメッセージを送った後で、ビットマップ内に新規の書き込み要求によって影響を受ける記憶装置の領域を前記第 1 のコントローラが蓄積する段階と、

前記第 2 の記憶装置が書き込み要求を処理することができるようになった後で、前記第 1 の記憶装置から前記第 2 の記憶装置へコピーされることになる前記第 1 の記憶装置の記憶領域を前記第 1 のコントローラが前記ビットマップを用いて識別する段階と、

前記第 1 の記憶装置の識別された領域の内容を前記第 2 のコントローラが前記第 2 の記憶装置にコピーする段階と、

を含む方法。

【請求項 16】

揮発性記憶装置及び不揮発性記憶装置を具備する第 1 の記憶装置と、

揮発性記憶装置及び不揮発性記憶装置を具備する第 2 の記憶装置と、

前記第 1 の記憶装置と関連付けられた第 1 のコントローラと、

前記第 2 の記憶装置と関連付けられた第 2 のコントローラと、

を備え、

前記第 1 のコントローラは、

前記第 1 の記憶装置において書き込み要求を受け取り、

前記第 1 の記憶装置の揮発性記憶装置又は不揮発性記憶装置に前記書き込み要求を保存するとともに、前記第 1 の記憶装置の揮発性記憶装置に保存されている書き込み要求を前記第 1 の記憶装置の不揮発性記憶装置に対して周期的にコミットすることによって、前記第 1 の記憶装置で受け取られた前記書き込み要求を処理し、

前記第 1 の記憶装置によって処理された書き込み要求を識別する基準ラベルに基づく情報である 1 つの書き込み要求を指定する情報と共に前記第 2 の記憶装置に対して同期コミットメッセージを周期的に送る、

ように構成されており、

前記第 2 のコントローラは、

前記第 2 の記憶装置において書き込み要求を受け取り、

前記第 2 のコントローラが、前記第 2 の記憶装置の揮発性記憶装置又は不揮発性記憶装置を用いて、前記同期コミットメッセージにおいて識別された前記書き込み要求を処理した後で、揮発性ディスク記憶装置において保存されている書き込み要求を前記第 2 の記憶装置の不揮発性ディスク記憶装置にコミットし、

前記同期コミットメッセージを受け取って該同期コミットメッセージ内の情報によって識別された前記書き込み要求を処理した後で、前記処理された書き込み要求が前記第 2 の記憶装置の不揮発性記憶装置に正常にコミットされた場合に、前記指定された書き込み要求に先行していたすべての書き込み要求に関連するデータが、前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことと前記指定された書き込み要求に関連するデータが前記第 2 の記憶装置の不揮発性記憶装置に書き込まれたことを示す情報を含む確認メッセージを送るよう構成されていることを特徴とするミラーデータ記憶システム。

【請求項 17】

前記第 1 のコントローラは、前記第 1 の記憶装置で処理された書き込み要求によって影響を受ける記憶装置の領域を識別するように更に構成されていることを特徴とする請求項 16 に記載のシステム。

【請求項 18】

前記第 1 のコントローラは、第 1 のビットマップ内に前記記憶装置の識別された領域を蓄積するように構成されていることを特徴とする請求項 17 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

本発明はディスクドライブ又は他の記憶装置のミラーコピーを生成するための手法に関する。

【背景技術】

【0002】

多くのコンピュータシステムでは、同一のデータを複数の記憶装置の各々に記憶することによってフォールト・トレランスのレベルが設定される。同一のデータを有する記憶装置はミラー装置と呼ばれ、ミラーセットに属すると言われる。ミラーセットの1つのミラー装置が障害を発生するか又はアクセス不能となった場合、ミラーセットの他の1つ又は複数のミラー装置が継続してデータにアクセスする。

【0003】

ミラーセットの各装置で同一のデータを保持するためには、各装置は、ミラーセットにデータを保存するために全ての要求(すなわち、全ての書き込み要求を)を受け取って処理しなければならない。ミラーセット内の装置は、該装置がこうした書き込み要求を処理できない場合には、ミラーセットの他の装置と不一致となる。ミラーセットのメンバーが不一致になると、1つのミラー装置から別のミラー装置にデータをコピーするためにミラーセットコピーを実行することができる。ミラーセットコピーを保持するための1つの手法においては、コンピュータシステムは停止され、1つのミラー装置から他のミラー装置に全てのデータがコピーされる。

【発明の開示】

【課題を解決するための手段】

【0004】

1つの一般的な態様において、コンピュータシステムにおいて第1の記憶装置のミラーコピーが第2の記憶装置に保持される。第1の記憶装置は、関連するコントローラを含み、第2の記憶装置は、関連するコントローラと、揮発性記憶装置と、不揮発性記憶装置とを含む。記憶装置で受け取られた書き込み要求が処理される。同期コミットメッセージが、書き込み要求を指定する情報と共に第2の記憶装置に送られ、第2の記憶装置のコントローラが、同期コミットメッセージを受け取った後に、指定された書き込み要求に関連するデータが第2の記憶装置の不揮発性記憶装置に書き込まれたことを確認する。

【0005】

実施形態は、1つ又はそれ以上の次の特徴を含む。例えば、第2の記憶装置のコントローラは、指定された書き込み要求に先行する全ての書き込み要求に関連するデータが第2の記憶装置の不揮発性記憶装置に書き込まれたことを確認することができる。或いは、第2の記憶装置のコントローラは、指定された書き込み要求を処理することができ、且つ指定された書き込み要求と先行する書き込み要求と関連するデータが第2の記憶装置の不揮発性記憶装置に書き込まれたことを確認することができる。第2の記憶装置のコントローラは、第2の記憶装置の揮発性記憶装置の正常なキャッシュフラッシュを確認することができる。

【0006】

同期コミットメッセージと共に送られた情報は、第1の記憶装置によって処理されたか、或いは処理されることになる書き込み要求を識別する基準ラベルとすることができる。基準ラベルは、他の書き込み要求に割り当てられた基準ラベルに対して順次に割り当てることができる。第2の記憶装置で受け取られた全ての書き込み要求は、同期コミットメッセージの基準ラベルによって識別された書き込み要求の処理の前に順次処理することができる。

【0007】

書き込み要求によって影響を受ける記憶装置の識別された領域は、例えば第1のビットマップに蓄積することができる。同期コミットメッセージを送った後は、記憶装置の新規に識別された領域は第2のビットマップに蓄積することができる。処理された書き込み要求のデータが第2の記憶装置の不揮発性記憶装置に書き込まれたことを、第2の記憶装置のコントローラが確認した後で、書き込みデータが不揮発性記憶装置に正常に書き込まれ

10

20

30

40

50

たことを示すために、ステータスメッセージを第1の記憶装置に送ることができる。書き込みデータが正常に書き込まれたことを示すステータスメッセージを受け取った後で、第1のビットマップを削除し、第2のビットマップを第1のビットマップとして指定することができる。

【0008】

第2の記憶装置が利用可能でなかった期間の後で、第1のビットマップの内容を、回復ビットマップにコピーし、次にこれを用いて第1の記憶装置から第2の記憶装置にコピーされることになる第1の記憶装置の記憶領域を識別することができる。第1の記憶装置の識別された記憶領域を第2の記憶装置にコピーすることができ、新規に受け取られた書き込み要求を第3のビットマップ内に第2の記憶装置において蓄積することができる。

10

【0009】

第2の記憶装置は、第1の記憶装置に対して上述したような1つ又はそれ以上の特徴及び機能を実行することができ、第1の記憶装置は第2の記憶装置に対して上述したような1つ又はそれ以上の特徴及び機能を実行することができる。

【0010】

別の一般的な態様において、コンピュータシステムにおける第1の記憶装置のミラーコピーを第2の記憶装置に保持することは、関連するコントローラと、揮発性記憶装置と、不揮発性記憶装置とを含む第1の記憶装置において書き込み要求を受け取る段階と、第1の記憶装置で受け取られた書き込み要求を処理する段階と、関連するコントローラと揮発性記憶装置と不揮発性記憶装置とを含む第2の記憶装置において書き込み要求を受け取る段階と、第2の記憶装置において受け取られた書き込み要求を処理する段階とを含む。第2の記憶装置が書き込み要求を処理することができなくなる期間に入ろうとしていることを第2の記憶装置が判定した後で、第1の記憶装置のコントローラは、書き込み要求を指定する情報と共に第2の記憶装置に対して同期コミットメッセージを送り、第2の記憶装置のコントローラは、同期コミットメッセージを受け取った後で、指定された書き込み要求に関連するデータが第2の記憶装置の不揮発性記憶装置に書き込まれたことを確認する。同期コミットメッセージを送った後で、第1の記憶装置のコントローラは、ビットマップ内に新規の書き込み要求によって影響を受ける記憶の領域を蓄積する。第2の記憶装置が書き込み要求を再度処理できるようになると、第1の記憶装置から第2の記憶装置へコピーされることになる第1の記憶装置の記憶の記憶領域を識別するように、第1の記憶装置のコントローラはビットマップを使用して、第1の記憶装置の識別された領域の内容を第2の記憶装置にコピーする。

20

30

【0011】

上述の手法の実施形態は、方法又はプロセス、装置又はシステム、或いはコンピュータがアクセス可能な媒体上のコンピュータソフトウェアを含むことができる。

【発明を実施するための最良の形態】

【0012】

1つ又はそれ以上の実施形態の詳細が添付図面並びに以下の説明で記載される。他の特徴及び利点は、以下の説明、図面、及び請求項から明らかになる。

各図面中、同一の参照記号は同一の要素を表す。

40

【0013】

図1は、第1のデータ記憶装置105と第2のデータ記憶装置110を含むミラーセット100のブロック図である。図1の実施形態において、データ記憶装置はディスクドライブである。他の実施形態では、データ記憶装置はディスクドライブのアレイ又は他の記憶装置とすることができる。

【0014】

説明を簡単にするため、ディスクの1つはマスタディスクとされ、主データ記憶装置として機能し、他のディスクはスレーブディスクとされ、冗長バックアップとして機能する。双方のディスクがアクティブで且つ同じデータを含む場合には、マスタ/スレーブのステータスは2つのディスクに対して任意に割り当てることができる。実際に以下に述べる

50

同期化技法に対して、2 - ディスクの実施形態では実際に2つのマスタ - スレーブ関係が維持され、各ディスクは一方の関係ではマスタとして、他方の関係ではスレーブとして機能する。図1においては、ディスク105はマスタディスクとして指定され、ディスク110はスレーブディスクとして指定されている。

【0015】

第1のI/O(「入力/出力」)コントローラ115は第1のディスク105と関連付けられ、第2のI/Oコントローラ120は第2のディスク110と関連付けられる。I/Oコントローラ115、120は、ディスク上のデータの読み込みと書き込みを制御する。

【0016】

例えば、プロセッサとすることができるクライアント125は、同一の書き込み要求130を両方のI/Oコントローラに送る。各書き込み要求はデータを含む。加えて、連続参照番号などの基準ラベルは、各書き込み要求に関連付けられる。I/Oコントローラは、書き込み要求からのデータをそれぞれのディスクに書き込むので、通常の条件下では両方のディスクは同一のデータを含む。典型的には、各I/Oコントローラは書き込み要求を同じ順番で処理する。これを達成するために、I/Oコントローラは書き込み要求を基準ラベルの順番に処理するので、これはI/Oコントローラが同一順番の書き込み要求を受け取る必要がないことを意味する。

【0017】

また、クライアント125は読み出し要求135をI/Oコントローラに送る。1つの実施形態において、両方のディスクが同一のデータを含む場合、マスタディスクだけが読み出し要求135に応答する。他の実施形態では、スレーブディスク又は両方のディスクが応答することができる。マスタディスクに障害が発生するか又はアクセス不能となった場合には、スレーブディスクがマスタディスクとして指定変更され、継続してデータをクライアント125に供給する。従って、ディスク105に障害が発生した場合にはディスク110がマスタディスクとなる。

【0018】

ミラーセット100のディスクは、ある時間期間の間に書き込み要求を処理することができない場合には、ピアのディスクとは不一致のデータを含むようになる。例えば、スレーブディスクがある時間期間の間に利用不能であった場合、スレーブディスクのデータはマスタディスクのデータとは異なるものとなる。ミラーセットのディスクが不一致になると、ミラーセットコピーを実行して、「正常な」データを有するディスクから不一致データを有するディスクにデータをコピーする。幾つかの大容量記憶装置においては、このプロセスは長時間を要する可能性があり、この間はミラーディスクが同一のデータを含まないことにより、システムのフォールト・トレランスのレベルが狭くなる。

【0019】

システムのフォールト・トレランスレベルを改善するために、「良好な」データを有するディスクから、不一致データを有するディスクに保存されていなかったデータ部分だけをコピーすることによって、両方のディスクが同一データを含む状態(回復と呼ぶ場合がある)にミラーディスクを復元するのに必要な時間を短縮することができる。ディスクの部分だけをコピーするこのプロセスは、増分不一致コピー又は差分コピー(ここで、差分とは、一方のディスクに加えられており、他のディスクにはない変更を意味する)と呼ぶことができる。

【0020】

一般に、増分不一致コピーは、1つ又はそれ以上のミラーディスクに対して行われるスレーブの変更によって行うことができ、その結果、利用不能期間の後で、「良好な」データだけを有するディスクから不一致データを有するミラーディスクに保存されていなかったデータをコピーすることによって、不一致データを有するミラーディスクを復元することができるようになる。一般に、ミラーディスクに対して加えられてきた変更の監視は、ミラーセット内の両方のミラーディスクが同一のデータを含むことが分かった時点(この

10

20

30

40

50

時点ではミラーディスクは同期していることができる)の後で加えられた変更の推移を把握することが必要とされる。

【0021】

同期化された時点の後にミラーディスクに加えられた変更の監視は、データが揮発性のディスクキャッシュには既書き込まれているが、ミラーディスクの不揮発性記憶装置には未だ書き込まれていない時に、システム、サブシステム、又はプロセッサが書き込み要求を完了したものと記録する際に問題となる可能性がある。2つ以上のRAID (Redundant Array of Inexpensive Disks) ディスクへの書き込みに要する時間が増大したことに起因して、書き込み要求がディスクキャッシュに出されてから、全てのデータの揮発性ディスク記憶装置への書き込みを完了するまでの時間期間はかなりのものとなる可能性がある。例えばRAID方式を使用して2つ以上のディスクにわたってミラーデータをストライピングした場合に、これは特に重大な問題となる恐れがある。

10

【0022】

ディスクキャッシュをフラッシュしてキャッシュ内のデータをディスク記憶装置にコミットすることにより、特定の書き込み要求の基準ラベルまでミラーディスク上のデータを周期的に同期化することで(例えば各ディスクが同一のデータを含む場合)増分不一致コピーの効果を改善することができる。ディスクキャッシュをフラッシュすることにより、処理されていた全ての書き込み要求が揮発性ディスク記憶装置内に確実に保存される。

20

【0023】

ミラーセット100内の各I/Oコントローラ115、120は、ビットマップ155、156内のディスクに対して加えられた変更を蓄積することによって、I/Oコントローラのそれぞれのディスクに対して行われた書き込み要求130の推移を把握する。ビットマップ155、156は、ディスクの各領域が書き込み要求150によって影響を受けたかどうかを表示するために1つ又はそれ以上のビットを使用するデータ構造である。この実施形態におけるビットマップはディスクに保存される。別の実施形態では、システムが停止するか、又はミラーセットに含まれていない揮発性記憶装置にビットマップを保存するまで、揮発性メモリにビットマップを保存することができる。ビットマップによって得られる抽象性(又は粒度)のレベルは、ビットで表わされる記憶領域のサイズに基づく。一般に、各ビットは、単一の書き込み要求によって記憶装置の対応する領域に書き込まれるデータよりも実質的に大きなデータを表す。ここでビットマップ155、156はディスク変更ビットマップと呼ばれることがある。

30

【0024】

ビットマップ及びディスクは、固有の識別子によって関連付けることができる。固有の識別子は、例えば、ビットマップとディスクが適用されるクライアント125のインスタンスを識別するディスク識別子を組み込むことができる。ビットマップとディスクを関連付けることにより、変更されたディスク領域を適切なディスクに確実にコピーすることができる。例えば、特定のビットマップと特定のディスク又はディスク上の特定のデータセットとの関連付けは、ミラーセットが取り外し可能ディスク(すなわち、コンピュータのハウジングユニットを開かずに取り外すことが可能なディスク)を含む場合に重要である。

40

【0025】

マスタI/Oコントローラと呼ぶことができる1つのI/Oコントローラが、スレーブI/Oコントローラと呼ぶことができる他のI/Oコントローラに同期コミットメッセージ160を周期的に送る。図示し且つ以下に説明されるように、第1のI/Oコントローラ115がマスタI/Oコントローラであり、第2のI/Oコントローラ120がスレーブI/Oコントローラである。しかしながら、ディスク110がマスタであり、ディスク105がスレーブという関係で、第2のI/Oコントローラ120が同時にマスタI/Oコントローラとして作用する(第1のI/Oコントローラが同時にスレーブI/Oコントローラとして作用する)点に留意することは重要である。

50

【 0 0 2 6 】

同期コミットメッセージ 1 6 0 は、ミラーディスク上のデータが同期化されるようになる書き込み要求の基準ラベルを識別する。第 1 の I / O コントローラ 1 1 5 は、ディスク変更ビットマップ 1 5 5 のバックアップコピー 1 6 5 を作成して同期化プロセスの間に障害が発生した場合に回復できるようにし、新しいディスク変更ビットマップを開始させて、次の同期化において使用されることになる後続の全ての書き込みを蓄積する。

【 0 0 2 7 】

第 2 の I / O コントローラ 1 2 0 が、第 1 の I / O コントローラ 1 1 5 によって送られた同期コミットメッセージを受け取ると、第 2 の I / O コントローラは、該第 2 の I / O コントローラが同期コミットメッセージにおいて識別された書き込み要求と全てのこれまでの書き込み要求とを既に処理したかどうかを判断する。処理していない場合、第 2 の I / O コントローラは、同期化を開始する前に、当該書き込み要求と全てのこれまでの書き込み要求とが処理されるまで待機する。

10

【 0 0 2 8 】

第 2 の I / O コントローラ 1 2 0 が、同期コミットメッセージで識別された書き込み要求と全てのこれまでの書き込み要求とを処理すると、又は同期コミットメッセージを受け取った時に、第 2 の I / O コントローラが書き込み要求と全てのこれまでの書き込み要求とを既に処理していた場合には、第 2 の I / O コントローラは、そのディスクコントローラキャッシュをフラッシュして処理された書き込み要求を不揮発性記憶装置にコミットする。キャッシュのフラッシュ化が正常である場合、第 2 の I / O コントローラ 1 2 0 は第 1 の I / O コントローラ 1 1 5 に確認メッセージ 1 7 0 を送る。フラッシュ及び同期化が正常であったことの確認を受け取るとすぐに、第 1 の I / O コントローラ 1 1 5 は、このディスク変更ビットマップのバックアップコピー 1 6 5 をクリアする。同期化が正常ではない場合、又は所定の時間内に第 1 の I / O コントローラ 1 1 5 が確認を受け取らない場合には、第 1 の I / O コントローラ 1 1 5 は、ビットマップ 1 5 5 とバックアップ 1 6 5 とを（通常は論理和を取ることによって）結合し、結合されたビットマップを第 2 のディスク 1 1 0 の復元に使用する。

20

【 0 0 2 9 】

特定の時点からディスクに加えられた変更だけを蓄積する増分不一致コピープロセスは、ディスク障害が検出されると開始する（従って、利用不能の期間の間に加えられた変更だけを蓄積する）ことができ、又は、システムがアクティブである時は何時でも用いることができる（従って、システム動作中に常時ミラーセットに加えられた変更を蓄積する）。利用不能の期間だけ変更が蓄積される時には、利用不能の期間は、利用可能になった後でディスクを回復する際に有効となる増分不一致コピープロセスが利用不能になっている（これは、「正常な停止」と呼ぶことができる、プロセス全体でディスクが利用不能になる場合に行われることがある）ディスクのディスクコントローラのキャッシュフラッシュから開始する必要がある。

30

【 0 0 3 0 】

別の実施形態は、異なるディスク変化ビットマップの同期化時点後に加えられたディスク変化の蓄積を開始するのではなく、ディスク変化ビットマップにおいて特定の書き込み要求を削除することによる変化を蓄積することを含むことができる。これは回復時間を短縮することができる。

40

【 0 0 3 1 】

ミラーセット 1 0 0 の I / O コントローラ 1 1 5 と 1 2 0 のいずれかは、同期コミットプロセスを開始して、両方のディスクが特定の書き込み要求の基準ラベル全体を通じて同一のデータを含むよう保証することができる。利用不能期間の後で、最後の同期化以降変更されたディスク領域だけをコピーすることによって、同一データを保存するミラーコピーであるように不一致データを有するミラーディスクを復元することができる。

【 0 0 3 2 】

1 つの実施形態では、日付（又は日付と時間）とミラーディスクの 1 つに対する特定の

50

データセットとの関連付けを含むことができる。これは、書き込み要求の基準ラベルが必ずしも固有のものでない場合に有益とすることができる。例えば、基準ラベルが、クライアントを制御するオペレーティングシステムのリセット時において、ある固定値（例えば1）から再起動する連続番号である場合、この書き込み要求の基準ラベルは固有ではない可能性がある。このような書き込み要求は、データセットを保存するディスクのディスクキャッシュがフラッシュされる時に、日付（又は日付と時間）をデータセットと関連付けることによって固有に識別することができる。別の方法として又は追加として、クライアントが再起動される時に、非固有の書き込み要求の基準ラベルの識別に役立つように、ミラーセットに新しいインスタンス番号を付与することができる。他の固有の識別子は、クライアント識別子、ミラーセット識別子、及びデータセット識別子を単独で又は組み合わせて含むことができる。

10

【0033】

例証として図1では、ミラー装置として2つのディスクを使用してミラーデータのセットを保存しているが、増分不一致コピーの利点は、この特定の実施形態に限定されるものではなく、RAID方式を含む別の数又は種類の記憶装置を伴う他の実施形態に対しても同様に適用可能である。例えば、別の実施形態は、3つ又はそれ以上のディスクをミラーリングすることができ、又はミラーディスクの複数のインスタンス化（例えば、4つのディスクを使用して同一のディスクに対して2つのミラーセットを提供することができる）を行うことができる。

【0034】

20

図2を参照すると、プロセス200は増分不一致追跡を使用して、ミラーディスクセットのディスクが、正常な停止を通じて利用不能になった期間中に、不一致になったミラーディスクセットへの同期化を復元する準備をする。図2で設定されたミラーディスクの実施形態は、別個のI/Oコントローラによって各々制御される2つのディスク記憶装置を有する。各I/Oコントローラは、プロセッサから同一の書き込み要求を受け取り、受け取った書き込み要求を起こった順番に処理する。基準ラベルは各書き込み要求と関連付けられ、書き込み要求の順序付けに使用される。別の実施形態では、特定の書き込み要求が処理されたか（完了したか）どうかを追跡することができる。これにより書き込み要求をバラバラの順序で（すなわち非順次で）処理することが可能となる。両方のディスクがアクティブである場合、各ディスクは同一のデータを含む。

30

【0035】

プロセス200は、ディスクの1つが利用不能の期間（ステップ205）に入っていると判定された時に開始される。この判定がなされると、利用不能になっているディスクのI/Oコントローラは、I/Oコントローラによって既に処理されている不揮発性記憶装置の書き込み要求に対してコミットするよう命令される（ステップ210）。利用不能となるディスクをスレーブディスク、アクティブなディスクをマスタディスクと呼び、これらに関連付けられたI/Oコントローラは、スレーブI/OコントローラとマスタI/Oコントローラと呼ぶことができる。マスタI/Oコントローラは、ディスク変更ビットマップ内のマスタディスクに加えられた変更の蓄積を開始し（ステップ220）、継続して書き込み要求を受け取って処理し（ステップ225）、ディスク変更ビットマップを更新し、処理された書き込み要求から結果として得られた全ての変更をマスタディスクに反映させる（ステップ230）。ディスク変更ビットマップの各ビットは、マスタディスクの領域を表す。他の実施形態では、ディスク変更ビットマップの各ビットによって表示されるディスクスペース量を変えることができる。

40

【0036】

マスタI/Oコントローラは又、継続してスレーブディスクのステータスを監視する（ステップ235）。スレーブディスクが使用可能となり、新規の書き込み要求の処理を開始すると、マスタI/Oコントローラは、図3に関連して以下に説明するように、回復プロセス300を開始する（ステップ240）。

【0037】

50

図3を参照すると、回復プロセス300は、ディスク変更ビットマップによって示されるように、マスタディスクのスレーブディスク部分へのコピーを含む。回復プロセスは、ミラーセットが継続して新規の書き込み要求を処理している間にアクティブなバックグラウンドプロセスとして行われる。回復プロセス300は、マスタI/Oコントローラがディスク変更ビットマップのバックアップコピーを作成し、ディスク変更ビットマップのオリジナルバージョンを回復ビットマップとして指定することで開始する(ステップ310)。また、マスタI/Oコントローラは、新しいディスク変更ビットマップを開始して、後続の全ての変更をマスタディスクに蓄積する(ステップ320)。ディスク変更ビットマップと新しいディスク変更ビットマップのバックアップコピーは、回復プロセスの間に障害が発生した場合の回復を可能にする。

10

【0038】

マスタI/Oコントローラは、回復ビットマップ内の各ビットを検査し(ステップ330)、対応するマスタディスクの領域が既に変更されたことをビットが示すかどうかを判定する(ステップ340)。示さない場合、マスタI/Oコントローラは回復ビットマップの次のビットの検査に進む(ステップ345)。マスタディスク領域が既に変更されたことをビットが示す場合、マスタI/Oコントローラは、後続の書き込み要求が、スレーブディスクの対応するディスク領域を変更したかどうかを判定する(ステップ345)。

【0039】

後続の書き込み要求が、対応するスレーブディスク領域を既に変更していた場合には、マスタI/Oコントローラは、後続の書き込み要求によって変更されていないスレーブディスク領域の部分に相当する、マスタディスク領域の部分だけをコピーする(ステップ350)。マスタI/Oコントローラは、スレーブI/Oコントローラに、回復プロセスの間にスレーブディスクにより処理される書き込み要求のリストを保持させることにより、変更されていない領域を識別することができ、リストの各エントリは修正された実メモリ部分を識別する。或いは、スレーブI/Oコントローラは、より微細な粒度を有するディスク変更ビットマップを保持することができるので、ビットマップの各々のビットは、書き込み要求の修正が許可されているディスクの最小の部分に対応するようにする。スペースを保護するために、スレーブI/Oコントローラは、粒度が変化するビットマップを保持して、ディスクの修正された部分に対してだけ微細な粒度のマップが保持されるようにすることができる。

20

30

【0040】

スレーブディスク領域に対して後続の変更が何も加えられなかった場合、マスタI/Oコントローラは、マスタディスクの全領域をスレーブディスクにコピーする(ステップ355)。

【0041】

マスタI/Oコントローラは、コピーされているデータの部分を修正し、後続の書き込み要求によって上書きされることになるデータ書き込みの潜在的な非効率さを排除する(ステップ345-355)。例えば、書き込み要求WR-102がディスク領域12に保存されたデータの部分を変更し、WR-155もまた、ディスク領域12の別の部分に保存されたデータを変更する場合には、ディスク領域12へのデータを書き込むプロセスは、各書き込み要求に必要とされる領域12の部分だけを変更することができる。

40

【0042】

追加又は代替として、スレーブI/Oコントローラは、コピーされているデータの部分を修正することができる。例えば、スレーブI/Oコントローラが、マスタディスクからコピーされたデータによって更新されることになる同一のディスク領域を修正する新規の書き込み要求を受け取った場合には、スレーブI/Oコントローラは、マスタディスクからコピーされているデータの部分を修正することができる。

【0043】

マスタディスク領域(又はその部分)をスレーブディスクにコピーした後、マスタI/Oコントローラは、回復ビットマップ内の別のビットを検査する必要があるかどうかを判

50

断し（ステップ360）、必要な場合には次のビットを検査する（ステップ330）。

【0044】

マスタI/Oコントローラが、回復ビットマップ内の全てのビットが検査されたことを判定した時に回復が完了する。終了すると、コピーされたデータをスレーブディスクにコミットするために、マスタI/Oコントローラは、任意選択的にスレーブディスク・キャッシュをフラッシュする同期化を開始することができる（ステップ370）。マスタI/Oコントローラが、後続の同期化又はフラッシュ化が正常でないと判定する（ステップ375）と、マスタI/Oコントローラはディスク変更ビットマップのバックアップコピーを新しいディスク変更ビットマップと結合し（通常は論理和を取ることによって）（ステップ380）結合されたディスク変更ビットマップを使用して回復プロセス300を繰り返す。スレーブI/Oコントローラの同期化とスレーブディスクコントローラのフラッシュ化が正常である場合には、マスタI/Oコントローラはバックアップディスク変更ビットマップをクリアする（ステップ390）。

10

【0045】

図2に関連して説明された実施形態は、マスタディスクからスレーブディスクにディスク領域をコピーする時に粒度のレベルを変更するが、別の実施形態では、後続の書き込み要求によって領域の部分が上書きされるかどうかに関係なく、常に変更された領域全体をコピーすることができる。図3の実施形態は、回復ビットマップを回復の間に2度使用できないように処理する。回復プロセスの間に障害からの回復を可能とするために、ビットを処理する前にディスク変更ビットマップのバックアップコピーを行う（ステップ310）。

別の実施形態では、回復ビットマップは回復の間に破棄されることはなく、回復ビットマップ自体を使用することによって、回復プロセス障害から回復することができる。この実施形態は、ビットを処理する前にディスク変更ビットマップのコピーを行うことはない（ステップ310）。別の代替実施形態は、スレーブディスクが使用可能に戻った後で、且つ回復プロセスが正常に完了する前に行われたマスタディスク変更を蓄積するために新しいディスク変更ビットマップを使用することはない。

20

【0046】

図4 - 図6を参照すると、増分不一致コピープロセスは、ミラーセットを使用する時は常にアクティブとすることができる。図4 - 図6のミラーセットの実施形態は、図2に関連して説明された方法で書き込み要求を受け取る2つのディスク記憶装置と2つのI/Oコントローラとを有する。

30

【0047】

各I/Oコントローラは、ディスク変更ビットマップ内のディスクに加えられた変更を蓄積することによって、I/Oコントローラのディスクに対して加えられた書き込み要求の推移を把握する。1つのI/Oコントローラ（マスタI/Oコントローラと呼ばれる）は、周期的に同期コミットメッセージを別のI/Oコントローラ（スレーブI/Oコントローラと呼ばれる）に送り、周期的な同期化プロセスを開始する。図4はマスタI/Oコントローラによって実行される周期的な同期化プロセスを示す。図5は、スレーブI/Oコントローラによって実行される周期的な同期化プロセスを示す。図6は、同一データを有するミラーコピーとなる、不一致データを有するミラーディスクを復元するためのプロセスを示す。

40

【0048】

図4を参照すると、マスタI/Oコントローラは、スレーブI/Oコントローラと周期的な同期化を行うためのプロセス400を開始する。プロセス400は、マスタI/Oコントローラが、プロセッサから書き込み要求を受け取って処理し（ステップ410）、ディスク変更ビットマップ内のマスタディスクに加えられた変更を蓄積する（ステップ415）ことで開始する。マスタI/Oコントローラは、ミラーセットを同期化すべきかどうかを判断する（ステップ420）。マスタI/Oコントローラは、例えば、最後の同期化から所定時間が経過した後、最後の同期化から所定の数の書き込み要求が処理された後、又は2つのミラーディスク間の増分不一致が一定割合となった後で、同期コミットメッセ

50

ージを作成することができる。同期化が要求される時を判定する際に、ミラーディスク間で同期化されていないデータ量に対して同期化周波数（これは、不揮発性記憶装置にディスクキャッシュが書き込まれている時間の間はディスクキャッシュのフラッシュ化が全ての書き込み要求の処理を停止するので、システムの性能を低下させる可能性がある）を平衡させることができる（これは、ミラーセットに対して同一のデータを復元するための増分不一致コピーを行うのにより長い時間を要する可能性がある）。

【0049】

ミラーセットを同期化すべきであることをマスタI/Oコントローラが判断すると、マスタI/Oコントローラは、同期コミットメッセージをスレーブI/Oコントローラに送る（ステップ430）。同期コミットメッセージは、ミラーディスク上のデータが同期化されることになる書き込み要求の基準ラベルを識別する。マスタI/Oコントローラは、ディスク変更ビットマップのバックアップコピーを作成（ステップ435）し、同期化プロセスで障害が発生した場合に回復できるようにし、新しいディスク変更ビットマップを開始し（ステップ440）て、次の同期化において使用するためにこの時点からマスタI/Oコントローラのディスクに加えられるディスク変更を蓄積する。マスタI/Oコントローラは、継続してプロセッサから書き込み要求を受け取って処理を行い（ステップ445）、マスタディスクの全ての変更を反映させるよう新しいディスク変更ビットマップを更新する（ステップ450）。

【0050】

スレーブI/Oコントローラによるキャッシュフラッシュと同期化とが正常であった（ステップ455）ことの確認を受け取ると、マスタI/Oコントローラは、バックアップディスク変更ビットマップをクリアし（ステップ460）、同期化のプロセスが終了する。

【0051】

或いは、例えば、マスタI/Oコントローラが所定の時間内にスレーブI/Oコントローラから確認メッセージを受け取らなかったか、又は同期化に障害が発生したというメッセージを受け取ったことにより、マスタI/Oコントローラは同期化が正常ではなかったことを判断することができる（ステップ455）。この場合には、マスタI/Oコントローラは、バックアップディスク変更ビットマップと新しいディスク変更ビットマップとを結合し（通常、論理和を取ることによって）（ステップ470）、スレーブI/Oコントローラとこれに関連するディスクとが動作可能であることを判断すると、結合されたディスク変更ビットマップを使用して、図3に関連して説明されたような回復プロセス300を開始し、マスタディスクからスレーブディスクにどのディスク領域をコピーすべきかを導く（ステップ475）。

【0052】

図5を参照すると、プロセス500は、ミラーディスク上のデータが同期化されることになる書き込み要求基準ラベルを識別する同期コミットメッセージを、スレーブI/Oコントローラが受け取った時に開始する（ステップ510）。スレーブI/Oコントローラは、同期コミットメッセージと全てのこれまでの書き込み要求とにおいて識別された書き込み要求を既に処理したかどうかを判定する（ステップ520）。処理していない場合には、スレーブI/Oコントローラは、同期化を開始する前に、該書き込み要求と全てのこれまでの書き込み要求が処理されるまで待機する。

【0053】

スレーブI/Oコントローラが、同期コミットメッセージで識別された書き込み要求と、全てのこれまでの書き込み要求とを処理すると、又は同期コミットメッセージが受け取られた時にスレーブI/Oコントローラが該書き込み要求と全てのこれまでの書き込み要求とを既に処理している場合には、スレーブディスクコントローラは、そのキャッシュをフラッシュして、処理された書き込み要求を不揮発性のディスク記憶装置にコミットし（ステップ530）、キャッシュフラッシュが正常であったかどうかを判定する（ステップ540）。キャッシュフラッシュが正常であった場合、スレーブI/Oコントローラは、

マスタI/Oコントローラに確認メッセージを送る(ステップ550)。キャッシュフラッシュが正常でなかった場合には、スレーブI/Oコントローラは、マスタI/Oコントローラに障害発生メッセージを送る(ステップ560)。マスタI/Oコントローラに適切なメッセージを送った後、スレーブI/Oコントローラはプロセス500を終了する。

【0054】

図6は、同一データを保存するミラーコピーであるように不一致データを有するミラーディスクを復元するためのプロセス600を示す。以下の説明は、ミラーセット内のディスクの1つ(スレーブディスク)がこれまでに障害が発生したか、或いは利用不能となっており、且つミラーセットが同期化された前回以降、他のアクティブなディスク(マスタディスク)に加えられた全ての変更を含むディスク変更ビットマップが存在することを仮定している。これは、例えば図4 - 図5に関連して説明されたプロセスを実行することによって達成することができる。

10

【0055】

スレーブディスクが利用可能でない場合、マスタI/Oコントローラは、継続してプロセッサから書き込み要求を受け取って処理し(ステップ610)、最後の同期化以降に加えられたディスク変更を追跡するディスク変更ビットマップのマスタディスクに対してに加えられた変更を蓄積する(ステップ620)。マスタI/Oコントローラが、スレーブディスクが回復して書き込み要求処理の開始が可能であると判断する(ステップ630)と、マスタI/Oコントローラは、ディスク変更ビットマップを使用して、図3に関連して説明されたような回復プロセス300を開始し、マスタディスクのデータと同一のデータを含むようにスレーブディスクを復元する(ステップ640)。

20

【0056】

プロセス400、500及び600を実行することによって達成された増分不一致コピーは、プロセス200を実行することによって達成されたものとは異なる。特にプロセス400 - 600は、ディスク変更ビットマップがミラーセットのアクティブな間に更新されるので、ディスクの1つに予期しないディスク障害又はコントローラ障害が発生している間にミラーディスクセットを再確立するのに有効である。プロセス200は、次のディスク利用不能期間の警告が、ディスクキャッシュのフラッシュ実行を可能にするのに十分であり、ディスク変更ビットマップに蓄積されるようになる残りのアクティブなディスクに対する変更が開始される際、ミラーディスクセットの再確立だけに有効である。しかしながら、プロセス200は特定の時間に実行されるだけであるので、オーバーヘッドの処理はプロセス400 - 600よりも大幅に少ない結果となる。

30

【0057】

図7を参照すると、回復プロセス700は、周期的な同期化の実行と、ディスク変更ビットマップによって示されるようなマスタディスクの部分のスレーブディスクへのコピーとを含む。回復プロセス700は、マスタI/Oコントローラがディスク変更ビットマップのバックアップコピーを作成して、ディスク変更ビットマップのオリジナルのバージョンを回復ビットマップとして指定(ステップ710)して開始する。マスタI/Oコントローラはまた、新しいディスク変更ビットマップを開始して、後続の全ての変更をマスタディスクに蓄積する(ステップ720)。

40

【0058】

図3に関連して説明されたように、マスタI/Oコントローラは、回復ビットマップの各ビットを検査して(ステップ730)、マスタディスク領域が変更されたことをビットが示す場合、マスタディスクは、該マスタディスクの変更された部分をスレーブディスクにコピーする(ステップ735)。

【0059】

図4に関連して上述したように、マスタI/OコントローラはスレーブI/Oコントローラとの同期化プロセスを周期的に開始する。特に、ミラーセットを同期化すべきであることをマスタI/Oコントローラが判断する(ステップ740)と、マスタI/Oコント

50

ローラは、同期コミットメッセージをスレーブI/Oコントローラに送り、ディスク変更ビットマップのコピーを作成し、新しいディスク変更ビットマップを開始して、この時点からマスタI/Oコントローラのディスクに加えられる変更を蓄積する(ステップ745)。

【0060】

スレーブディスクコントローラによるキャッシュフラッシュと同期化とが正常であった(ステップ750)という確認を受け取ると、マスタI/Oコントローラは、スレーブディスクに正常にコピーされていたマスタディスクの領域を示すビットを、回復ビットマップから除外する(ステップ755)。マスタI/Oコントローラは、例えば回復プロセスの間にマスタディスクによって処理されたビットのリストを保持し、リストされたビットを回復ビットマップから削除することによって、これを達成することができる。しかしながら、マスタI/Oコントローラが、同期化が正常でなかったと判定する(ステップ750)と、マスタI/Oコントローラは、バックアップディスク変更ビットマップと新しいディスク変更ビットマップとを結合し(ステップ760)、結合されたディスク変更ビットマップを使用して、図3に関連して説明された回復プロセス300を開始する(ステップ765)。

10

【0061】

マスタI/Oコントローラが、回復ビットマップの全てのビットを検査したと判定する(ステップ770)と、回復プロセスを終了し、バックアップディスク変更のビットマップをクリアする(ステップ775)。実施形態は、方法又はプロセス、装置又はシステム、若しくはコンピュータ媒体上のコンピュータソフトウェアを含むことができる。添付の請求項の精神及び範囲から逸脱することなく、種々の変更が可能であることは理解されるであろう。例えば、開示された手法のステップを別の順序で実行する場合、及び/又は開示されたシステムの構成要素を別の方法で結合する場合、及び/又は別の構成要素によって置き換えられるか補足される場合には、有利な結果が得られるであろう。

20

【図面の簡単な説明】

【0062】

【図1】ミラードライブシステムのブロック図である。

【図2】ミラーディスク間の差異を監視するプロセスのフローチャートである。

【図3】不一致になったミラーディスクのセットに対する同期性を回復するためのプロセスのフローチャートである。

30

【図4】マスタ入出力コントローラによって実行される周期的な同期化プロセスを示すフローチャートである。

【図5】スレーブ入出力コントローラによって実行される周期的な同期化プロセスを示すフローチャートである。

【図6】同一のデータを有するミラーコピーであるように、不一致データを有するミラーディスクを復元するためのプロセスのフローチャートである。

【図7】回復プロセスの間に実行される周期的な同期化を示すフローチャートである。

【符号の説明】

【0063】

- 115 I/Oコントローラ1
- 120 I/Oコントローラ2
- 125 クライアント
- 155 ビットマップ
- 156 ビットマップ
- 160 同期コミットメッセージ
- 165 バックアップビットマップ
- 170 確認メッセージ

40

【図1】

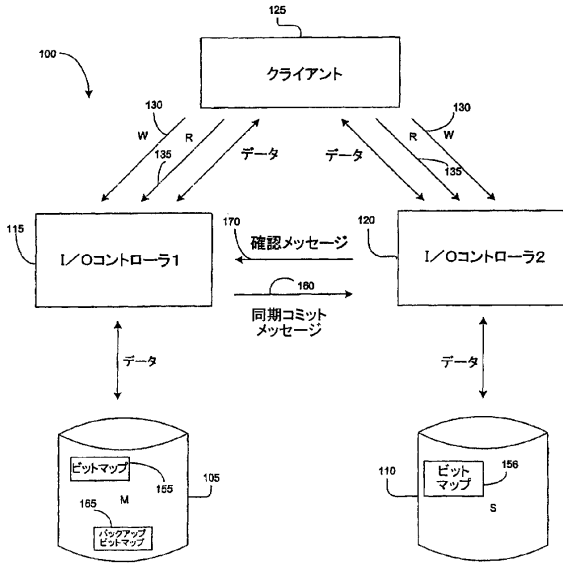


FIG. 1

【図2】

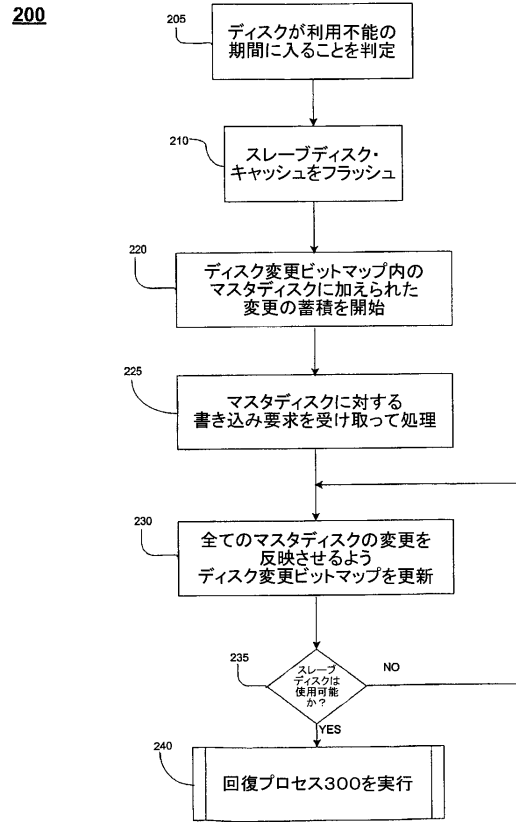


FIG. 2

【図3】

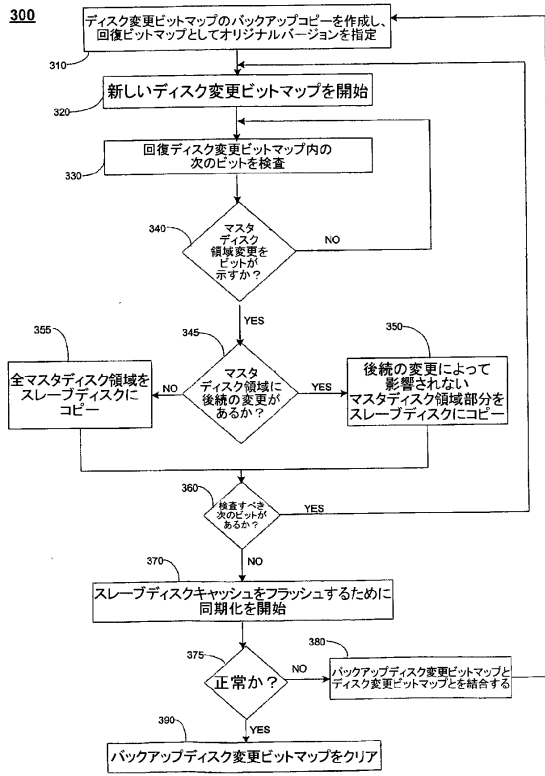


FIG. 3

【図4】

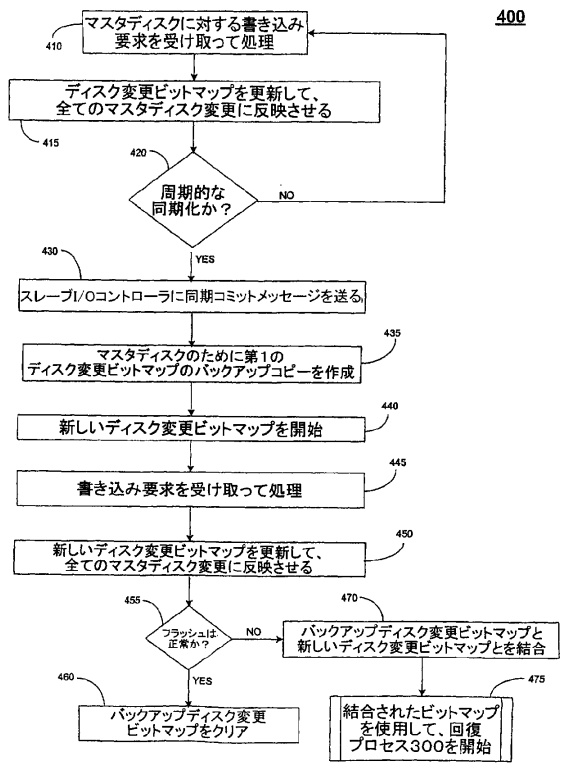


FIG. 4

【図5】

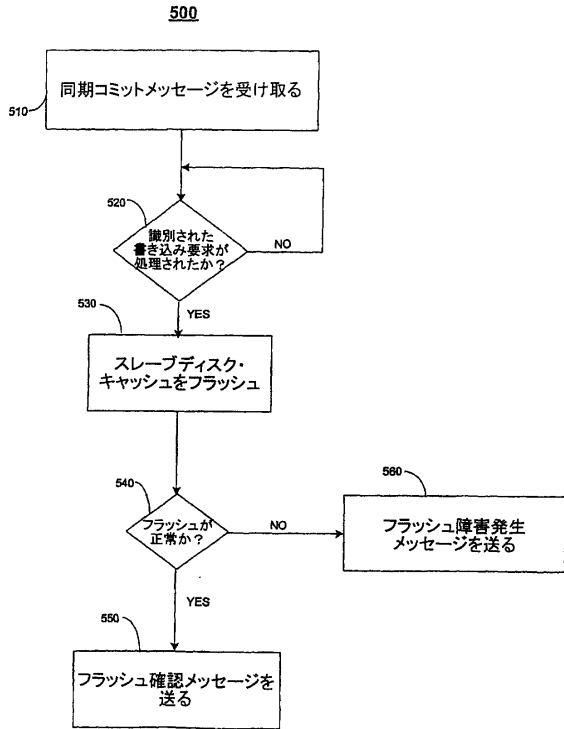


FIG. 5

【図6】

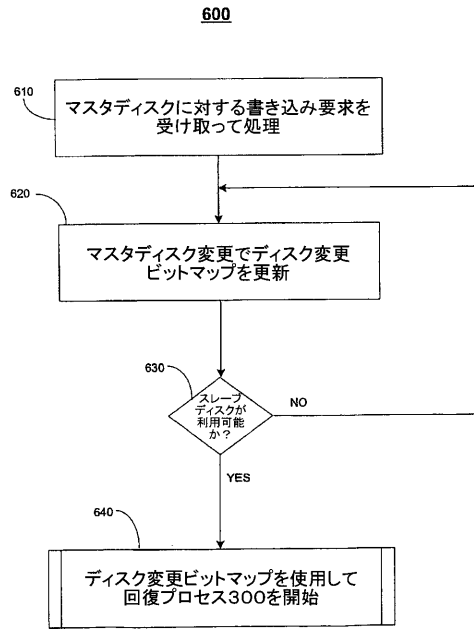


FIG. 6

【図7】

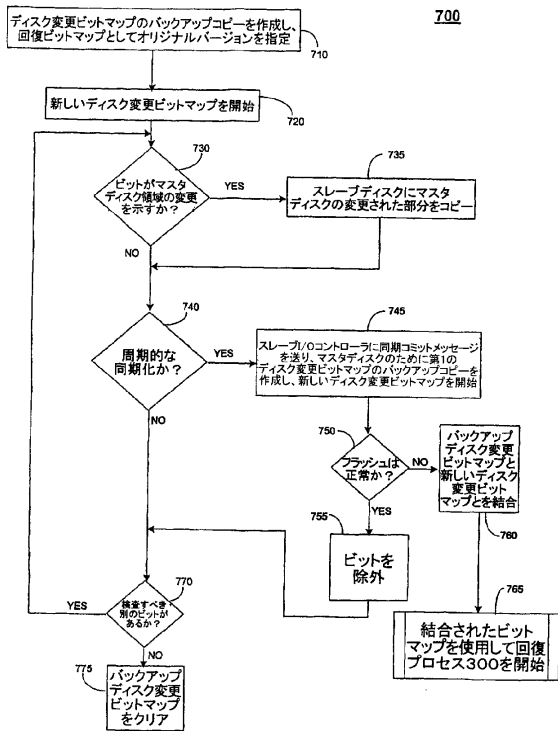


FIG. 7

フロントページの続き

- (72)発明者 トレンブレイ グレン エイ
アメリカ合衆国 マサチューセッツ州 01568 アップトン サウス ストリート 139
- (72)発明者 リヴィール ポール エイ
アメリカ合衆国 マサチューセッツ州 01519 グラフトン ストラットン ロード 12
- (72)発明者 カマン チャールズ エイチ
アメリカ合衆国 マサチューセッツ州 01773 リンカーン オーク メドー ロード 10
- (72)発明者 グラナン ゲアリー
アメリカ合衆国 マサチューセッツ州 01719 ボックスボロー レオナード ロード 35

審査官 上嶋 裕樹

- (56)参考文献 特表2001-501002(JP, A)
特開平2-166509(JP, A)
国際公開第01/040952(WO, A1)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

G06F 3/06