

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5147941号  
(P5147941)

(45) 発行日 平成25年2月20日 (2013. 2. 20)

(24) 登録日 平成24年12月7日 (2012. 12. 7)

(51) Int. Cl.

F I

G 0 6 F 3/06 (2006. 01)

G 0 6 F 3/06 3 0 4 F

G 0 6 F 13/10 (2006. 01)

G 0 6 F 13/10 3 4 0 A

G 0 6 F 12/00 (2006. 01)

G 0 6 F 12/00 5 3 1 D

請求項の数 9 (全 17 頁)

(21) 出願番号	特願2010-516502 (P2010-516502)	(73) 特許権者	390009531
(86) (22) 出願日	平成20年7月16日 (2008. 7. 16)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公表番号	特表2010-533911 (P2010-533911A)		I N T E R N A T I O N A L B U S I N E S S M A C H I N E S C O R P O R A T I O N
(43) 公表日	平成22年10月28日 (2010. 10. 28)		アメリカ合衆国 1 0 5 0 4 ニューヨーク州 アーモンク ニュー オーチャードロード
(86) 国際出願番号	PCT/EP2008/059305		
(87) 国際公開番号	W02009/010532		
(87) 国際公開日	平成21年1月22日 (2009. 1. 22)	(74) 代理人	100108501
審査請求日	平成23年6月8日 (2011. 6. 8)		弁理士 上野 剛史
(31) 優先権主張番号	11/780, 454	(74) 代理人	100112690
(32) 優先日	平成19年7月19日 (2007. 7. 19)		弁理士 太佐 種一
(33) 優先権主張国	米国 (US)	(74) 代理人	100091568
早期審査対象出願			弁理士 市位 嘉宏
			最終頁に続く

(54) 【発明の名称】 異なるネットワークを介した1次ストレージから2次ストレージへの書き込みコピーを管理するための方法、システム、およびコンピュータ・プログラム

(57) 【特許請求の範囲】

【請求項 1】

少なくとも1つの1次ストレージ、前記少なくとも1つの1次ストレージへの入力/出力 (I/O) アクセスを管理する少なくとも1つの1次デバイス、前記少なくとも1つの1次ストレージへの書き込みがコピーされる少なくとも1つの対応する2次ストレージ、及び前記少なくとも1つの対応する2次ストレージに対する入力/出力 (I/O) アクセスを管理する少なくとも1つの2次デバイスの第1グループと、少なくとも1つの1次ストレージ、前記少なくとも1つの1次ストレージへの入力/出力 (I/O) アクセスを管理する少なくとも1つの1次デバイス、前記少なくとも1つの1次ストレージへの書き込みがコピーされる少なくとも1つの対応する2次ストレージ、及び前記少なくとも1つの対応する2次ストレージに対する入力/出力 (I/O) アクセスを管理する少なくとも1つの2次デバイスの、第2グループと、に関する情報を維持することを含み、

前記維持される情報は、前記第1グループの前記少なくとも1つの1次デバイスが前記第1グループの前記2次デバイスと通信するように第1のネットワーク・プロトコルを使用して動作可能であり、前記第2グループの前記少なくとも1つの1次デバイスが前記第2グループの前記2次デバイスと通信するように第2のネットワーク・プロトコルを使用して動作可能であることを示し、

前記1次デバイスから前記2次デバイスへの書き込みをコピーするために、障害の前記1次または2次グループ内の前記1次デバイスのうちの1つから障害通知を受信すること

前記維持される情報から前記第1のネットワーク・プロトコルを識別すること、

前記維持される情報から前記第1のネットワーク・プロトコルを識別することに応答して、第1のネットワーク・プロトコルを使用して、第1のネットワークを介し、前記第1のグループ内の前記少なくとも1つの1次デバイスへ、前記対応する少なくとも1つの2次デバイスへの書き込みコピーを停止するためのフリーズ・コマンドを発行すること、および

前記維持される情報から前記第2のネットワーク・プロトコルを識別することに応答して、第2のネットワーク・プロトコルを使用して、第2のネットワークを介し、前記第2のグループ内の前記少なくとも1つの1次デバイスへ、前記対応する少なくとも1つの2次デバイスへの書き込みコピーを停止するためのフリーズ・コマンドを発行することを含み、

10

前記第1および第2のネットワーク・プロトコルを使用して、前記第1および第2のグループ内の前記1次デバイスから、前記フリーズ・コマンドが受信された旨の肯定応答を受信すること、

前記1次デバイスで書き込みを完了するため、および、前記フリーズ・コマンドの送信先であった前記第1および第2のグループ内のすべての前記デバイスからの肯定応答の受信に応答して、変更記録データ構造内に前記完了された書き込みを示すための信号を、前記1次デバイスに送るために、前記第1のネットワーク・プロトコルを使用して、前記第1のグループ内の前記少なくとも1つの1次デバイスに、実行コマンドを発行すること、

前記1次デバイスで書き込みを完了するため、および、前記フリーズ・コマンドの送信先であった前記第1および第2のグループ内のすべての前記デバイスからの肯定応答の受信に応答して、変更記録データ構造内に前記完了された書き込みを示すための信号を、前記1次デバイスに送るために、前記第1のネットワーク・プロトコルを使用して、前記第2のグループ内の前記少なくとも1つの1次デバイスに、実行コマンドを発行すること、および

20

前記障害通知をもたらした障害からの回復に応答して、前記変更記録データ構造内に示された書き込みをそれらの対応する2次デバイスにコピーするために、それぞれ前記第1および第2のネットワーク・プロトコルを使用して、前記第1および第2のグループ内の前記1次デバイスに信号を送ること、

をさらに含む、

30

方法。

#### 【請求項2】

前記第1および第2のグループ内の前記1次ストレージへの従属的書き込みが、前記第1および第2のグループ内のいずれかの前記1次ストレージに従属的書き込みが書き込まれた順に、前記対応する2次ストレージにコピーされるように、前記第1および第2のグループ内の前記1次ストレージを、少なくとも1つの整合性グループ内で維持すること、をさらに含む、請求項1に記載の方法。

#### 【請求項3】

前記第1および第2のグループ内の前記1次ストレージへの前記書き込みが、前記対応する2次ストレージでの前記書き込みが完了した旨の肯定応答が受信されるまで、完了しないように、前記データが、前記第1および第2のグループ内の前記1次ストレージならびに前記対応する2次ストレージに、同期的に書き込まれる、請求項2に記載の方法。

40

#### 【請求項4】

前記第1のグループおよび第2のグループ内の前記デバイスが異機種のデバイスを備える、請求項1からの請求項3までのいずれかに記載の方法。

#### 【請求項5】

少なくとも1つの1次デバイスと、該少なくとも1つの1次デバイスによって管理される少なくとも1つの対応する1次ストレージからなる第1のグループにおける前記1次デバイスが前記1次ストレージに対する書き込みが対応する2次ストレージにコピーできないと判断することに応答して、前記1次デバイスから、制御システムに、第1のネットワ

50

ーク・プロトコルを使用して、障害通知を通信することと、

少なくとも1つの1次デバイスと、該少なくとも1つの1次デバイスによって管理される少なくとも1つの対応する1次ストレージからなる第2のグループにおける前記1次デバイスが前記1次ストレージに対する書き込みが対応する2次ストレージにコピーできないと判断することに応答して、前記1次デバイスから、制御システムに、第2のネットワーク・プロトコルを使用して、障害通知を通信することと、

前記制御システムから、前記第1のネットワーク・プロトコルを使用して、前記第1のグループにおける前記少なくとも1つの1次デバイスにおいて、フリーズ・コマンドを受信することと、

前記制御システムから、前記第2のネットワーク・プロトコルを使用して、前記第2のグループにおける前記少なくとも1つの1次デバイスにおいて、フリーズ・コマンドを受信することと、

前記フリーズ・コマンドの受信に応答して、前記第1及び第2グループにおける前記1次デバイスから前記対応する2次ストレージに対する書き込みのコピーを停止することを含み、

前記フリーズ・コマンドの受信に応答して、前記第1のネットワーク・プロトコルを使用し、前記第1のネットワークを介して前記フリーズ・コマンドが受信された旨の肯定応答を、前記制御システムに送信すること、

前記フリーズ・コマンドの受信に応答して、前記第2のネットワーク・プロトコルを使用し、前記第2のネットワークを介して前記フリーズ・コマンドが受信された旨の肯定応答を、前記制御システムに送信すること、

前記第1および第2のグループ内のすべての前記1次デバイスから前記フリーズ・コマンドが受信された旨の前記肯定応答を前記制御システムが受信するのに応答して、前記第1のネットワーク・プロトコルを使用して前記制御システムから実行コマンドを受信すること、

前記実行コマンドの受信に応答して、前記第1および第2のグループ内の前記対応する1次ストレージへの書き込みを完了すること、および

前記書き込みの完了に応答して、変更記録データ構造内に完了された1つの書き込みを示すこと、

をさらに含む、

方法。

#### 【請求項6】

前記第1及び第2のグループにおける前記1次デバイスにおける任意のものに対する依存する書き込みの順序が、前記対応する2次ストレージに対する依存する書き込みのコピーにおいて維持されるように、前記第1及び第2のグループにおける前記1次ストレージに対する書き込みをコピーするステップをさらに有する、

請求項5に記載の方法。

#### 【請求項7】

前記第1及び第2のグループにおける前記1次ストレージに対する書き込みが、書き込みが前記対応する2次ストレージで完了したという確認が得られるまで完了しないように、前記第1及び第2のグループにおける前記1次デバイスから前記対応する2次ストレージに同期してデータがコピーされる、請求項6に記載の方法。

#### 【請求項8】

前記対応する2次ストレージへの前記書き込みコピーの完了に応答して、前記第1のネットワークを介し、前記第1のネットワーク・プロトコルを使用して、前記制御システムに書き込み完了メッセージを送信すること、および

前記対応する2次ストレージへの前記書き込みコピーの完了に応答して、前記第2のネットワークを介し、前記第2のネットワーク・プロトコルを使用して、前記制御システムに書き込み完了メッセージを送信すること、

をさらに含む、請求項5から請求項7のいずれかに記載の方法。

10

20

30

40

50

## 【請求項 9】

請求項 1 から請求項 8 までの任意の請求項の方法を実行するように適合されたプログラム・コードを含み、コンピュータ上で動作する、コンピュータ・プログラム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、異なるネットワークを介した 1 次ストレージから 2 次ストレージへの書き込みコピーの管理のための、方法、システム、および製品 (an article of manufacture) に関する。

## 【背景技術】

10

## 【0002】

災害回復システムは、通常、単一のポイント・イン・タイム (point-in-time) での突発性破局障害、または経時的なデータ損失という、2 つのタイプの障害に対処する。第 2 のタイプの漸次的災害では、ボリュームに対する更新が失われる可能性がある。データ更新の回復を支援するために、リモート位置でデータのコピーを提供することができる。こうした 2 重コピーまたはシャドー・コピーは、通常、アプリケーション・システムが新しいデータを 1 次ストレージ・デバイスに書き込む際に作成される。2 次サイトでデータのリモート・コピーを維持するために、インターナショナル・ビジネス・マシーンズ社 (「IBM(R)」) の Extended Remote Copy (XRC)、Coupled XRC (CXRC)、Global Copy、および Global Mirror Copy などの、様々なコピー技術を使用することができる。

20

## 【0003】

データ・ミラーリング・システムでは、データはボリューム・ペア内に維持される。ボリューム・ペアは、1 次ストレージ・デバイス内のボリューム、および、1 次ボリューム内に維持されるデータの同一コピーを含む 2 次ストレージ・デバイス内の対応ボリュームからなる。1 次および 2 次のストレージ・コントローラを使用して、1 次および 2 次のストレージ・デバイスへのアクセスを制御することができる。

## 【0004】

データベース・システムなどの多くのアプリケーション・プログラムでは、ある種の書き込みは、前の書き込みが発生しない限り発生不可能であり、そうでなければデータ保全性が危険にさらされることになる。その保全性が前のデータ書き込みの発生に依存している、こうしたデータ書き込みは、従属的書き込み (dependent write) として知られている。1 次および 2 次のストレージ内のボリュームは、すべての書き込みがそれらの論理順に転送された場合、すなわち、すべての従属的書き込みがそれに依存する書き込みよりも先に転送された場合、整合性を持つ。整合性グループとは、従属的書き込みが整合的に保護されるような 1 次ボリュームへの更新の集まりである。整合時刻とは、2 次ボリュームに対する更新が整合性を持つことをシステムが保証する、最終時刻である。整合性グループは、1 次デバイスに書き込まれた順にリモートまたは 2 次サイトに書き込まれたポイント・イン・タイム時点での、すべての従属的書き込みを含む。さらに整合性グループは、整合性タイム・スタンプと同じかまたはこれより早期のタイム・スタンプを有する、整合性グループ内のすべてのデータ書き込みに関する整合時刻を有する。整合性グループは、ボリュームおよびストレージ・デバイス全体にわたって、データの整合性を維持する。したがって、データが 2 次ボリュームから回復した場合、回復したデータは、整合性グループのポイント・イン・タイム時点での整合性を持つことになる。

30

40

## 【0005】

整合性グループは、セッション内で形成される。セッションに割り当てられたすべてのボリューム・ペアは、同じ整合性グループ内に維持されるそれらの更新を有することになる。したがって、整合性グループ内でグループにまとめられることになるボリュームを決定するために、セッションが使用される。ジャーナルから、整合性グループから集められた更新が、2 次ボリュームに適用される。回復動作中、ジャーナルからの更新が 2 次ボリ

50

ュームに適用されている間にシステムが障害を起こした場合、2次ボリュームへの書き込みを完了しなかった更新は、ジャーナルから回復し、2次ボリュームに適用することが可能である。

【発明の概要】

【発明が解決しようとする課題】

【0006】

異なるネットワークを介した1次ストレージから2次ストレージへの書き込みコピーを管理するための方法、システム、および製品を提供する。

【課題を解決するための手段】

【0007】

異なるネットワークを介した1次ストレージから2次ストレージへの書き込みコピーを管理するための方法、システム、および製品、ならびにコンピュータ・プログラムが提供される。少なくとも1つの1次ストレージ、少なくとも1つの1次ストレージへの入力/出力(I/O)アクセスを管理する少なくとも1つの1次デバイス、および、少なくとも1つの1次ストレージへの書き込みがコピーされる少なくとも1つの対応する2次ストレージの、第1グループと、少なくとも1つの1次ストレージ、少なくとも1つの1次ストレージへの入力/出力(I/O)アクセスを管理する少なくとも1つの1次デバイス、および、少なくとも1つの1次ストレージへの書き込みがコピーされる少なくとも1つの対応する2次ストレージの、第2グループとで、情報が維持される。1次デバイスから2次デバイスへの書き込みをコピーするために、障害の1次または2次グループ内の1次デバイスのうちの1つから障害通知が受信される。第1のネットワーク・プロトコルを使用して、第1のネットワークを介し、第1のグループ内の少なくとも1つの1次デバイスへ、対応する少なくとも1つの2次デバイスへの書き込みコピーを停止するためのフリーズ・コマンドが発行される。第2のネットワーク・プロトコルを使用して、第2のネットワークを介し、第2のグループ内の少なくとも1つの1次デバイスへ、対応する少なくとも1つの2次デバイスへの書き込みコピーを停止するためのフリーズ・コマンドが発行される。

【0008】

他の実施形態では、第1および第2のグループ内の1次ストレージへの従属的書き込みが、第1および第2のグループ内のいずれかの1次ストレージに従属的書き込みが書き込まれた順に、対応する2次ストレージにコピーされるように、第1および第2のグループ内の1次ストレージが、少なくとも1つの整合性グループ内で維持される。

【0009】

他の実施形態では、第1および第2のグループ内の1次ストレージへの書き込みが、対応する2次ストレージでの書き込みが完了した旨の肯定応答が受信されるまで、完了しないように、データは、第1および第2のグループ内の1次ストレージならびに対応する2次ストレージに、同期的に書き込まれる。

【0010】

他の実施形態では、書き込みをログ記録するための要求が、第1または第2のネットワークを介して1次デバイスのうちの1つから受信される。書き込みは、書き込みログに記録される。第1および第2のグループ内のいずれかの1次デバイスに関して、いずれかの保留中のログ記録された書き込みが存在するかどうかに関する判別が実行される。それに関するログ完了が戻された書き込みよりも早期のポイント・イン・タイムを有する、第1および第2のグループ内の1次デバイスに関する書き込みログ内に、保留中のログ記録された書き込みが存在しない旨の決定に応答して、1次デバイスが対応する第2のデバイスへ書き込みをコピーできるようにするために、第1または第2のネットワークを介してログ完了が1次デバイスに戻される。対応する2次デバイスに書き込みがコピーされた旨の通知の受信に応答して、1つのログ記録された書き込みが書き込みログから除去される。

【0011】

他の実施形態では、第1および第2のネットワーク・プロトコルを使用して、第1およ

10

20

30

40

50

び第2のグループ内の1次デバイスから、フリーズ・コマンドが受信された旨の肯定応答が受信される。1次デバイスで書き込みを完了するため、および、フリーズ・コマンドの送信先であった第1および第2のグループ内のすべてのデバイスからの肯定応答の受信に  
10 応答して、変更記録データ構造内に完了された書き込みを示すための信号を、1次デバイスに送るために、第1のネットワーク・プロトコルを使用して、第1のグループ内の少なくとも1つの1次デバイスに、実行コマンドが発行される。1次デバイスで書き込みを完了するため、および、フリーズ・コマンドの送信先であった第1および第2のグループ内のすべてのデバイスからの肯定応答の受信に  
20 応答して、変更記録データ構造内に完了された書き込みを示すための信号を、1次デバイスに送るために、第1のネットワーク・プロトコルを使用して、第2のグループ内の少なくとも1つの1次デバイスに、実行コマンドが発行される。障害通知をもたらした障害からの回復に  
30 応答して、変更記録データ構造内に示された書き込みをそれらの対応する2次デバイスにコピーするために、それぞれ第1および第2のネットワーク・プロトコルを使用して、第1および第2のグループ内の1次デバイスに信号が送られる。

【0012】

他の実施形態では、第1のグループおよび第2のグループ内のデバイスは異機種のデバイスを備える。

【0013】

異なるネットワークを介した1次ストレージから2次ストレージへの書き込みをコピーするための方法、システム、および製品、ならびにコンピュータ・プログラムが提供される。1次ストレージへの書き込みが対応する2次ストレージへコピーできない旨の1次  
20 デバイスの決定に応答して、第1のネットワーク・プロトコルを使用して、少なくとも1つの1次デバイスと少なくとも1つの1次デバイスによって管理される少なくとも1つの対応する1次ストレージとの第1のグループ内の1次デバイスから、制御システムへと、障害通知が送られる。対応する1次ストレージへの書き込みが対応する2次ストレージへ  
30 コピーできない旨の1次デバイスの決定に応答して、第2のネットワーク・プロトコルを使用して、少なくとも1つの1次デバイスと少なくとも1つの1次デバイスによって管理される少なくとも1つの対応する1次ストレージとの第2のグループ内の1次デバイスから、制御システムへと、障害通知が送られる。第1のグループ内の少なくとも1つの1次  
40 デバイスで、第1のネットワーク・プロトコルを使用して、制御システムからフリーズ・コマンドが受信される。第2のグループ内の少なくとも1つの1次デバイスで、第2のネットワーク・プロトコルを使用して、制御システムからフリーズ・コマンドが受信される。第1および第2のグループ内の1次ストレージから対応する2次ストレージへの書き込み  
50 のコピーは、フリーズ・コマンドの受信に応答して中断される。

【0014】

他の実施形態では、第1および第2のグループ内のいずれかの1次ストレージへの従属的書き込みの順序が、対応する2次ストレージへの従属的書き込みのコピーにおいて保持されるように、第1および第2のグループ内の1次ストレージへの書き込みがコピーされる。

【0015】

他の実施形態では、第1および第2のグループ内の1次ストレージへの書き込みが、対応する2次ストレージでの書き込みが完了した旨の肯定応答が受信されるまで、完了しないように、データは、第1および第2のグループ内の1次ストレージから対応する2次  
40 ストレージへと、同期的にコピーされる。

【0016】

他の実施形態では、第1のグループ内の1次デバイスのうちの1つによって、書き込み要求が受信される。第1のネットワークを介し、第1のネットワーク・プロトコルを使用して、制御システムへの書き込みをログ記録するようにメッセージが送信され、ここで第1のグループ内の1次  
50 デバイスは、第1のネットワークを介して書き込みがログ記録された旨の制御システムからの肯定応答を受信するまで、対応する2次ストレージへの書き込

みをコピーしない。第2のネットワークを介し、第2のネットワーク・プロトコルを使用して、制御システムへの書き込みをログ記録するようにメッセージが送信され、ここで第2のグループ内の1次デバイスは、第2のネットワークを介して書き込みがログ記録された旨の制御システムからの肯定応答を受信するまで、対応する2次ストレージへの書き込みをコピーしない。第1および第2のグループ内の1次デバイスは、第1および第2のグループ内の1次デバイスによって、早期ポイント・イン・タイムを有する書き込みがそれらの対応する2次ストレージにコピーされるまでは、それらの対応する2次ストレージへの従属的書き込みをコピーしない。

【0017】

他の実施形態では、対応する2次ストレージへの書き込みコピーの完了に応答して、第1のネットワークを介し、第1のネットワーク・プロトコルを使用して、制御システムに書き込み完了メッセージが送信される。対応する2次ストレージへの書き込みコピーの完了に応答して、第2のネットワークを介し、第2のネットワーク・プロトコルを使用して、制御システムに書き込み完了メッセージが送信される。

【0018】

他の実施形態では、フリーズ・コマンドの受信に応答して、第1のネットワーク・プロトコルを使用し、第1のネットワークを介してフリーズ・コマンドが受信された旨の肯定応答が、制御システムに送信される。フリーズ・コマンドの受信に応答して、第2のネットワーク・プロトコルを使用し、第2のネットワークを介してフリーズ・コマンドが受信された旨の肯定応答が、制御システムに送信される。第1および第2のグループ内のすべての1次デバイスからフリーズ・コマンドが受信された旨の肯定応答を制御システムが受信するのに応答して、第1のネットワーク・プロトコルを使用して制御システムから実行コマンドが受信される。第1および第2のグループ内のすべての1次デバイスからフリーズ・コマンドが受信された旨の肯定応答を制御システムが受信するのに応答して、第2のネットワーク・プロトコルを使用して制御システムから実行コマンドが受信される。実行コマンドの受信に応答して、第1および第2のグループ内の対応する1次ストレージへの書き込みが完了される。書き込みの完了に応答して、変更記録データ構造内に完了された1つの書き込みが表示される。

【0019】

他の実施形態では、第1のグループ内の少なくとも1つの1次デバイスおよび1次ストレージは、第2のグループ内の少なくとも1つの1次デバイスおよび少なくとも1つの1次ストレージに関して異機種のデバイスを備え、ここで第1および第2のグループは、障害通知の送信動作、フリーズ・コマンドの受信動作、および書き込みコピーの中断動作を実行するための、異機種のストレージ・マネージャ・プログラムを有する。

【0020】

次に、本発明の好ましい諸実施形態について、以下の図面を参照しながら単なる例として説明する。

【図面の簡単な説明】

【0021】

【図1】ネットワーク・コンピューティング環境の実施形態を示す図である。

【図2】整合性グループ・メンバ情報の実施形態を示す図である。

【図3】書き込みログ・エントリ情報の実施形態を示す図である。

【図4】書き込み要求を処理するための動作の実施形態を示す図である。

【図5】書き込みが完了した旨の肯定応答を処理するための動作の実施形態を示す図である。

【図6】1つの2次デバイスの可用性において、障害を処理するための動作の実施形態を示す図である。

【図7】フリーズ・コマンド受信の肯定応答を処理するための動作の実施形態を示す図である。

【図8】諸実施形態の特定の説明された態様が内部に実装されるコンピュータ・アーキテ

10

20

30

40

50

クチャを示すブロック図である。

【発明を実施するための形態】

【0022】

図1は、ネットワーク・コンピューティング環境の実施形態を示す。1つまたは複数の1次デバイス2の第1のグループは、それぞれ1次ストレージ4への入力/出力(I/O)アクセスを管理し、2次デバイス6は、それぞれ2次ストレージ8へのI/Oアクセスを管理する。各1次デバイス2は、結合された1次ストレージ4への書き込みを、対応する2次デバイス6の2次ストレージ8に格納するために、対応する2次デバイス6にミラーリングする。1次デバイス2の第1のグループおよび対応する2次デバイス6は、第1のネットワーク・プロトコルを使用し、第1のネットワーク10を介して通信する。1つまたは複数の1次デバイス12の第2のグループは、それぞれ1つまたは複数の1次ボリューム16を有する1次ストレージ14への入力/出力(I/O)アクセスを管理し、2次デバイス18は、それぞれ1つまたは複数の2次ボリューム22を有する2次ストレージ20へのI/Oアクセスを管理する。第2のグループ内の各1次デバイス12は、含まれる結合されたボリューム16への書き込みを、対応する2次デバイス18の対応する2次ボリューム22に格納するために、対応する2次デバイス18にミラーリングする。1次デバイス12の第2のグループおよび対応する2次デバイス18は、第2のネットワーク・プロトコルを使用し、第2のネットワーク24を介して通信する。

10

【0023】

ネットワーク10および24の両方に結合された制御システム26は、異なるネットワーク10および24内にある1次ストレージ4および1次ボリューム16のうちのいずれかへのいずれかの書き込みが、あるポイント・イン・タイム時点で整合性を持つように、単一の整合性グループ内の1次ストレージ4および1次ボリューム16を管理する、制御ソフトウェア28を含む。このようにして、1次ストレージ4または1次ボリューム16への従属的書き込みは、それらが1次サイトへ書き込まれた順に、それらの対応する2次ストレージ8または2次ボリューム22へとミラーリングされる。整合性グループ内の第1および第2のグループ内のいずれの1次デバイス2および12における、その後の従属的書き込みも、整合性グループ内のいずれの1次デバイス2および12における以前の書き込みが完了するまでは、対応する2次デバイス6および18に書き込まれない。制御ソフトウェア26は、第1および第2の両方のネットワーク・プロトコルを使用して、それぞれ第1および第2のネットワーク10および24上で通信することができる。

20

30

【0024】

制御ソフトウェア28は、1つの整合性グループに含まれる1次ボリューム16/2次ボリューム22および1次ストレージ4/2次ストレージ8のあらゆるペアに関する情報を有する、整合性グループ情報30を維持する。さらに制御ソフトウェア28は、書き込みログ32内で保留されている、1次ボリューム16および1次ストレージ4への書き込みに関する情報もログ記録する。一実施形態では、1次デバイス2および12は、対応する2次ストレージ8、20に書き込みが正常にミラーリングされた旨を1次デバイス2および12が確認するまで、書き込みが完了しないように、それらの1次ストレージ4および1次ボリューム16へ同期的にデータを書き込む。

40

【0025】

ネットワーク10および24は、ストレージ・エリア・ネットワーク(SAN)、ローカル・エリア・ネットワーク(LAN)、イントラネット、インターネット、ワイド・エリア・ネットワーク(WAN)、ピアツーピア・ネットワーク、無線ネットワーク、アービトラート型(arbitrated)ループ・ネットワークなどを、含むことができる。説明された諸実施形態では、異なるネットワーク通信プロトコルを使用して、第1のネットワーク10および第2のネットワーク24上で通信する。たとえば一実施形態では、イーサネットおよびTCP/IPなどのパケットまたはステートレス(stateless)通信プロトコルを使用して第1のネットワーク10上で通信し、ストレージ・デバイス通信プロトコルを使用してファイバ・チャネル、シリアル接続SCSI(SAS)などの第2のネットワー

50



ク 2 4 上で通信することができる。

【 0 0 2 6 】

1 次デバイス 2、1 2 および 2 次デバイス 6、1 8 は、それぞれ、オペレーティング・システム 3 4、3 6、3 8、および 4 0 を含む。1 次デバイス 2 の第 1 のグループおよびそれらの対応する 2 次デバイス 8 は、制御ソフトウェア 2 8 と通信するため、および、1 次ストレージ 4 への書き込み要求と 2 次ストレージ 8 への書き込みのミラーリングとを管理するために、それぞれ、ストレージ・デバイス・ドライバ 4 2 および 4 4 を含む。1 次デバイス 1 2 の第 2 のグループおよびそれらの対応する 2 次デバイス 1 8 は、制御ソフトウェア 2 8 と通信するため、および、1 次ストレージ 1 4 への書き込み要求と 2 次ストレージ 2 0 への書き込みのミラーリングとを管理するために、それぞれ、ストレージ・マネージャ 4 6 および 4 8 を含む。1 次デバイス・ドライバ 4 2 およびストレージ・マネージャ 4 6 は、F R E E Z E / R U N モードで動作中の場合などの、対応する 2 次デバイス 6 および 1 8 への接続が使用不可の場合、完了した 1 次ストレージ 4 およびボリューム 1 6 への書き込みを示すように、変更記録ビットマップ 5 0 および 5 2 を維持する。

10

【 0 0 2 7 】

ストレージ 4、8、1 4、および 2 0 は、ハード・ディスク・ドライブ、フラッシュ・メモリなどの、単一のストレージ・デバイス、あるいは、J u s t a B u n c h o f D i s k s ( J B O D )、ネットワーク接続ストレージ ( N A S )、ハード・ディスク・ドライブ、直接アクセス・ストレージ・デバイス ( D A S D )、R e d u n d a n t A r r a y o f I n d e p e n d e n t D i s k s ( R A I D ) アレイ、仮想化デバイス、テープ・ストレージ、フラッシュ・メモリなどの、ストレージ・デバイスのアレイを、含むことができる。1 次デバイス 2 および 1 2 は、単一システム内に実装された複数の論理区画 ( L P A R ) または仮想プロセッサを備えることができる。

20

【 0 0 2 8 】

一実施形態では、第 1 のグループ内の 1 次デバイス 2 および対応する 2 次デバイス 6 がサーバを備え、ストレージ 4 および 8 は、内部または外部バス、シリアル・インターフェース、ユニバーサル・シリアル・バス ( U S B )、ファイアワイヤ・インターフェースなどを介して、デバイス 2、6 に接続された、デバイス 2、6 に対してローカルなハード・ディスク・ドライブを備えることができる。別の方法として、第 1 のグループ内のデバイス 2、6 およびストレージ 4、8 の組み合わせは、ネットワーク接続ストレージ ( N A S ) を備えることができる。一実施形態では、第 2 のグループ内のデバイス 1 2 および 1 8 は、R A I D アレイ、J B O D などの複数の論理ボリューム 1 6 および 2 2 を実装する、相互接続されたストレージ・デバイスを備える、ストレージ・システム 1 4 および 2 0 へのアクセスを管理する、エンタープライズ・ストレージ・サーバを備えることができる。

30

【 0 0 2 9 】

さらに、一実施形態では、第 1 のグループ内のデバイス 2、6 と共に使用される、1 つまたは複数のオペレーティング・システム 3 4、3 8、あるいはストレージ 4、8、またはそれらすべてが、第 2 のグループ内のデバイス 1 2 および 1 8 と共に使用される、オペレーティング・システム 3 6、4 0、あるいはストレージ 1 4、2 0、またはそれらすべてに関して、異機種である。一実施形態では、データの書き込みおよびミラーリングを管理するために使用されるストレージ・マネージャ・コードが、接続されたストレージ 4、8 に関するデバイス・ドライバ 4 2 内に実装される。一実施形態では、ストレージ・マネージャ・コード 4 6 および 4 8 は、エンタープライズ・ストレージ・サーバ内で使用されるハードウェアおよびソフトウェアの組み合わせを備えることができる。

40

【 0 0 3 0 】

図 2 は、整合性グループ内で管理される 1 次 / 2 次のストレージ・ペアに関する整合性グループ情報 3 0 内の、整合性グループ・メンバ・エントリ 7 0 に含めることが可能な、情報例を示す。エントリ 7 0 は、整合性グループ内の 1 次ストレージ 7 4 へのアクセスを管理する 1 次デバイス 7 2 と、1 次ストレージ 7 4 への書き込みのミラーリング先である対応する 2 次ストレージ 7 8 へのアクセスを管理する 2 次デバイス 7 6 と、1 次デバイス

50

72のネットワーク・アドレス80と、1次デバイス72と通信するために使用されるネットワーク・プロトコル82とを含む。

【0031】

図3は、1次ストレージ94への書き込みを実行する1次デバイス92、および書き込みのポイント・イン・タイム96を含む、書き込みログ32内の書き込みログ・エントリ90に含めることが可能な、情報例を示す。

【0032】

図4は、書き込み要求を処理するために、第1および第2のグループ内の1次デバイス2、12内のデバイス・ドライバ42およびストレージ・マネージャ46、ならびに制御ソフトウェア28によって実行される、動作の実施形態を示す。書き込み要求を受信すると(ブロック100)、デバイス・ドライバ42/ストレージ・マネージャ46は、制御ソフトウェア28への書き込みをログ記録するために、1次デバイスによって使用されるネットワーク10または24を介してメッセージを送信する(ブロック102)。書き込みをログ記録するためのメッセージの受信(ブロック104)にตอบสนองして、制御ソフトウェア28は、要求された書き込みに関してログ・エントリ90(図3)を書き込みログ32に追加する(ブロック106)。早期のポイント・イン・タイム96を有する第1および第2のグループ内の任意の1次デバイスに関して、保留中のログ記録された書き込みがない場合(ブロック108)、制御ソフトウェア28はログ完了を戻す(ブロック110)。さもなければ、完了していない早期のポイント・イン・タイムを有する保留中の書き込みが存在する場合(ブロック108)、制御は完了を戻さずに終了するため、結果として、1次デバイス2は、時間的に早期に、それらのそれぞれの2次ストレージ8またはボリューム22に書き込みがコピーされるまで、書き込みをコピーすることができない。ログ完了の受信にตอบสนองして、デバイス・ドライバ42/ストレージ・マネージャ46は、ログ記録された書き込みを、対応する2次ストレージへのアクセスを管理する2次デバイス6を介して、1次ストレージ4および対応する2次ストレージ8にコピーすることができる(ブロック112)。

【0033】

ある実施形態では、1次ストレージ・デバイス4は、対応する2次ストレージ8にデータが正常にコピーされるまで書き込みが完了しないように、データを同期的に書き込むことができる。一実施形態では、書き込みデータが2次ストレージ8内に格納されるまで、書き込みは完了しない。代替実施形態では、書き込みデータが、2次ストレージ8に書き込まれる前に、対応する2次ストレージ8へのアクセスを管理する2次デバイス6のキャッシュ内に格納された場合、書き込みを完了することができる。

【0034】

図5は、2次ストレージ8への書き込みの完了を処理するために、第1および第2のグループ内の1次デバイス2、12内のデバイス・ドライバ42/ストレージ・マネージャ46、ならびに制御ソフトウェア28によって実行される、動作の実施形態を示す。デバイス・ドライバ42/ストレージ・マネージャ46が対応する2次デバイスへの書き込みのコピーが完了した旨の肯定応答を受信すると(ブロック150)、1次デバイスによって使用されるネットワーク10、24を介して、書き込みが完了した旨のメッセージを制御ソフトウェア28に送信する(ブロック152)。書き込みが完了した旨の肯定応答の受信(ブロック154)にตอบสนองして、制御ソフトウェア28は、完了した書き込みに関するログ・エントリ90を書き込みログ32から削除する(ブロック156)。制御ソフトウェア28は、最も早期のポイント・イン・タイム96(図3)を有する、書き込みログ32内にログ記録された書き込み90を決定し(ブロック158)、決定されたログ記録された書き込みの1次デバイス2、12に、ログ完了を戻し(ブロック160)、フィールド92に示されたように、1次デバイス2、12が2次デバイス6、18への書き込みをコピーできるようにする。これらの動作によって、第1および第2のグループ内の1次ストレージ4およびボリューム16のいずれかへのその後の書き込みが、それらの対応する2次ストレージ8または2次ボリューム22に順序外れでコピーされないことが保証さ

10

20

30

40

50

れる。代替実施形態では、拡張された長期ビジー期間の使用など、様々な技法を使用して、データが順序外れで書き込まれないように保証することが可能であり、その結果、他の１次デバイス時間にそれらの早期書き込みを完了できるように、１次デバイスがある拡張された長期ビジー期間だけ書き込みのコピーを遅延させる。

【００３５】

図６は、１次デバイス２、１２が２次デバイス６、１８と通信できないことを示す障害通知を処理するために、第１および第２のグループ内の１次デバイス２、１２内のデバイス・ドライバ４２／ストレージ・マネージャ４６、ならびに制御ソフトウェア２８によって実行される、動作の実施形態を示す。１次デバイス２、１２が２次デバイス６、１８と通信する機能における障害を検出する（ブロック２００）か、または１次デバイス２ハードウェアによって通知されると、デバイス・ドライバ４２／ストレージ・マネージャ４６は、１次デバイス２、１２によって使用されるネットワーク１０、２４を介して、制御ソフトウェア２８に障害通知を送信する（ブロック２０２）。この障害は、１次デバイス２、１２と２次デバイス６、１８との間のネットワーク接続における障害、あるいは、たとえば２次デバイス６、１８または２次ストレージ８、２０の障害などの、２次サイトでの障害の結果である可能性がある。第１のネットワーク１０および第２のネットワーク２４を介して、第１および第２のグループ内の１次デバイス２、１２のうちのいずれかからの可能性がある障害通知を受信すると（ブロック２０４）、制御ソフトウェア２８は、第１のネットワーク・プロトコルを使用して、２次ストレージ８、２０への書き込みコピーを停止するように、第１のグループ内のそれぞれの１次デバイス２にフリーズ・コマンドを発行する（ブロック２０６）。さらに制御ソフトウェア２８は、第２のネットワーク・プロトコルを使用して、対応する少なくとも１つの２次デバイスへの書き込みコピーを停止するように、第２のグループ内のそれぞれの１次デバイス１２にフリーズ・コマンドを発行する（ブロック２０８）。制御ソフトウェア２８は、整合性グループ情報３０内のエントリ７０から、使用するための１次デバイス７２（図２）およびネットワーク・プロトコル８２を決定することができる。

【００３６】

フリーズ・コマンドを受信すると（ブロック２１０）、デバイス・ドライバ４２／ストレージ・マネージャ４６は、対応する２次ストレージ８または２次ボリューム２２への書き込みのコピーを中断する（ブロック２１２）。デバイス・ドライバ４２／ストレージ・マネージャ４６は、フリーズ・コマンドが受信された旨の肯定応答を、フリーズ・コマンドを送信した１次デバイス２、１２に送信する（ブロック２１４）。

【００３７】

図７に関して、制御ソフトウェア２８が、１次デバイス２、１２のうちの１つからフリーズ・コマンドの肯定応答を受信すると（ブロック２５０）、制御ソフトウェア２８は、第１のネットワーク１０および第２のネットワーク２４を介して、整合性グループ内のすべての１次デバイス２、１２から肯定応答が受信されたかどうかを判別する（ブロック２５２）。受信されなかった場合、制御は終了するか、または、すべての１次デバイス２、１２から肯定応答が受信されなかった場合、他の適切な処置を取ることができる。さもなければ、すべての１次デバイス２、１２から肯定応答が受信された場合、制御ソフトウェア２８は、第１のネットワーク・プロトコルを使用して、対応する２次デバイス６、１８への書き込みコピーを停止するように、第１のグループ内のそれぞれの１次デバイス２に実行コマンドを発行する（ブロック２５４）。制御ソフトウェア２８は、第２のネットワーク・プロトコルを使用して、対応する２次デバイス１８への書き込みコピーを停止するように、第２のグループ内のそれぞれの１次デバイス１２に実行コマンドを発行する（ブロック２５６）。制御ソフトウェア２８は、整合性グループ情報３０から実行コマンドに関する１次デバイスを決定することができる。

【００３８】

実行コマンドを受信すると（ブロック２５８）、デバイス・ドライバ４２／ストレージ・マネージャ４６は、１次ストレージ４または１次ボリューム１６への書き込みを完了し

10

20

30

40

50

(ブロック260)、変更記録データ構造50、52内に完了された書き込みを示す(ブロック262)。2次デバイス6、18あるいはストレージ8、20またはその両方が回復した後、1次デバイス2、12は、1次ストレージ4、14および2次ストレージ8、20を同期させるために、それらを介して変更記録ビットマップ50、52に示された書き込みをコピーすることができる。

#### 【0039】

説明された諸実施形態は、異なるネットワーク通信プロトコルを使用する異なるネットワークを介して配布された2次ストレージまたはボリュームへの、1次ストレージまたはボリュームへの書き込みのコピーを管理する。さらに説明された諸実施形態は、異なるネットワーク・プロトコルを使用して、異なるネットワーク内で1次デバイスを管理することによって、2次サイトでの障害を処理する。他の諸実施形態は、異なるネットワーク内の2次ストレージで、データがあるポイント・イン・タイム時点で整合性を持つように維持する。

#### 【0040】

追加の諸実施形態の詳細

説明された動作は、ソフトウェア、ファームウェア、またはそれらの任意の組み合わせを生成するために、標準的なプログラミングあるいはエンジニアリング、またはその両方の技法を使用する、方法、装置、または製品として実装可能である。説明された動作は、「コンピュータ読み取り可能媒体」内に維持されるコードとして実装可能であり、プロセッサは、コンピュータ読み取り可能媒体からコードを読み取って実行することができる。コンピュータ読み取り可能媒体は、磁気ストレージ媒体(たとえば、ハード・ディスク・ドライブ、フレキシブル・ディスク、テープなど)、光ストレージ(CD-ROM、DVD、光ディスクなど)、揮発性および不揮発性のメモリ・ディスク(たとえば、EEPROM、ROM、PROM、RAM、DRAM、SRAM、フラッシュ・メモリ、ファームウェア、プログラマブル論理など)などの、媒体を含むことができる。さらに、説明された動作を実装するコードは、ハードウェア論理(たとえば、集積回路チップ、プログラマブル・ゲート・アレイ(PGA)、特定用途向け集積回路(ASIC)など)内に実装することが可能である。さらに、説明された動作を実装するコードは、「伝送信号」内に実装することが可能であり、伝送信号は、大気を通じて、または光ファイバ、銅線などの伝送媒体を介して、伝搬可能である。内部にコードまたは論理が符号化される伝送信号は、無線信号、衛星伝送、電波、赤外線信号、Bluetooth(R)などを、さらに含むことができる。内部にコードまたは論理が符号化される伝送信号は、送信局による送信および受信局による受信が可能であり、伝送信号内部に符号化されたコードまたは論理を復号し、受信および送信の局またはデバイス側のハードウェアまたはコンピュータ読み取り可能媒体に格納することができる。「製品」は、内部にコードを実装可能な、コンピュータ読み取り可能媒体、ハードウェア論理、あるいは伝送信号、またはそれらすべてを含むことができる。動作の説明された諸実施形態を実装するコードが内部に符号化されたデバイスは、コンピュータ読み取り可能媒体またはハードウェア論理を含むことができる。もちろん、当業者であれば、本発明の範囲を逸脱することなく、この構成に対して多くの修正が実行可能であること、および、製品が当分野で知られた好適な情報担持媒体を含むことが可能であることを、理解されよう。

#### 【0041】

「実施形態」、「諸実施形態」、「1つまたは複数の実施形態」、「いくつかの実施形態」、および「一実施形態」という用語は、特に明示的に指定されていない限り、「本発明の(すべてではないが)1つまたは複数の実施形態」を意味する。

#### 【0042】

「含む」、「備える」、「有する」(「including」、「comprising」、「having」)およびそれらの変形用語は、特に明示的に指定されていない限り、「含むが限定されない」ことを意味する。

#### 【0043】

列挙されたアイテムのリストは、特に明示的に指定されていない限り、アイテムのいずれかまたはすべてが相互に排他的であることを示唆するものではない。

【 0 0 4 4 】

「ある」および「その」（「 a 」、 「 a n 」、および「 t h e 」）という用語は、特に明示的に指定されていない限り、「 1 つまたは複数」を意味する。

【 0 0 4 5 】

互いに通信しているデバイスは、特に明示的に指定されていない限り、互いに継続的に通信している必要はない。加えて、互いに通信しているデバイスは、直接、あるいは、 1 つまたは複数の媒介を介して間接的に、通信することができる。

【 0 0 4 6 】

互いに通信しているいくつかのコンポーネントを備える実施形態の説明は、こうしたコンポーネントのすべてが必要であることを示唆するものではない。これに対して、様々なオプションのコンポーネントは、本発明の多種多様な可能な実施形態を示すために記述されている。

【 0 0 4 7 】

さらに、プロセス・ステップ、方法ステップ、アルゴリズムなどは一連の順序で説明可能な場合があるが、こうしたプロセス、方法、およびアルゴリズムは、代替の順序で作動するように構成することが可能である。言い換えれば、説明可能なステップの任意の配列または順序は、必ずしもそのステップがその順序で実行される要件を示すものではない。本明細書で説明されるプロセスのステップは、事実上、任意の順序で実行可能である。さらに、いくつかのステップを同時に実行可能である。

【 0 0 4 8 】

本明細書において単一のデバイスまたは製品が説明される場合、単一のデバイス / 製品の代わりに、複数のデバイス / 製品が（それらが協働しているか否かに関わらず）使用可能であることが容易に理解されよう。同様に、本明細書において複数のデバイスまたは製品が（それらが協働しているか否かに関わらず）説明される場合、複数のデバイスまたは製品の代わりに単一のデバイス / 製品が使用可能であること、あるいは、示された数のデバイスまたはプログラムの代わりに異なる数のデバイス / 製品が使用可能であることが、容易に理解されよう。別の方法として、デバイスの機能あるいは特徴またはその両方は、こうした機能 / 特徴を有するものとして明示的に説明されていない 1 つまたは複数の他のデバイスによって、具体化可能である。したがって、本発明の他の実施形態は、デバイスそれ自体を含む必要はない。

【 0 0 4 9 】

図 4、5、6、および 7 に示された動作は、ある順序で発生するある種のイベントを示す。代替の諸実施形態では、ある種の動作は、異なる順序での実行、修正、または除去が可能である。さらに諸ステップを前述の論理に追加し、説明された諸実施形態に依然として合致することが可能である。さらに、本明細書で説明された動作は順番に発生可能であるか、またはある種の動作は並行して処理することが可能である。さらに動作は、単一の処理ユニットによって、または分散型処理ユニットによって、実行可能である。

【 0 0 5 0 】

図 8 は、デバイス 2、6、12、18、および 26（図 1）内で全体または部分的に実装可能な、コンピューティング・システム・アーキテクチャ 300 の実施形態を示す。アーキテクチャ 300 は、1 つまたは複数のプロセッサ 302（たとえば、マイクロプロセッサ）、メモリ 304（たとえば揮発性メモリ・デバイス）、およびストレージ 306（たとえば、磁気ディスク・ドライブ、光ディスク・ドライブ、テープ・ドライブなどの、不揮発性ストレージ）を含むことができる。ストレージ 306 は、内部ストレージ・デバイス、あるいは接続型ストレージまたはネットワーク・アクセス可能ストレージを含むことができる。ストレージ 306 内のプログラムは、当分野で知られた様式で、メモリ 304 にロードされプロセッサ 302 によって実行される。アーキテクチャは、ネットワークを介した通信を実行可能にするための 1 つまたは複数のアダプタ 308 をさらに含む。入

10

20

30

40

50

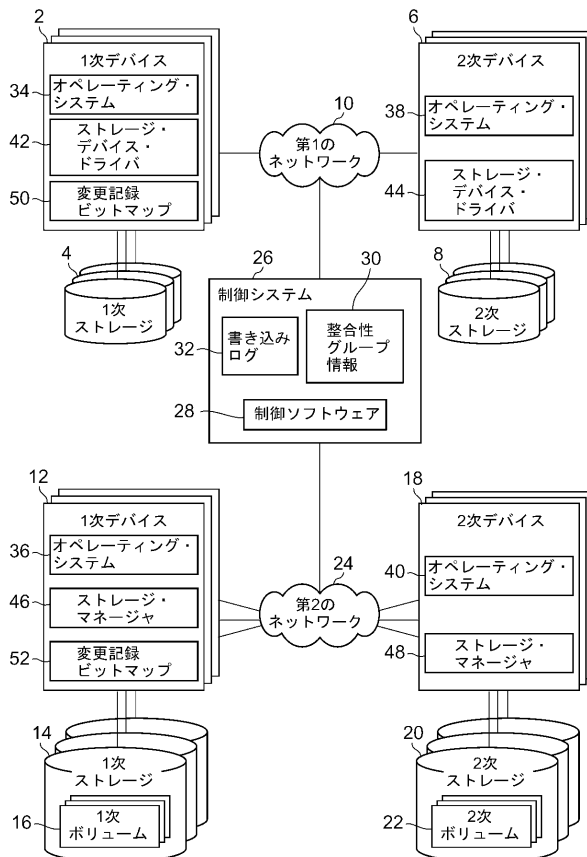
カデバイス 310 は、プロセッサ 302 へユーザ入力を提供するために使用することが可能であり、キーボード、マウス、ペン・スタイラス、マイクロフォン、タッチ・センシティブ・ディスプレイ・スクリーン、あるいは、当分野で知られた任意の他の起動または入力メカニズムを含むことが可能である。出力デバイス 312 は、プロセッサ 302、または、ディスプレイ・モニタ、プリンタ、ストレージなどの他のコンポーネントから送信される情報を、レンダリングすることができる。

#### 【0051】

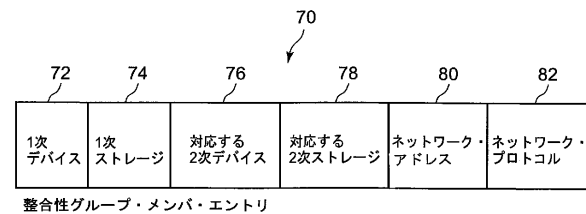
本発明の様々な実施形態の前述の説明は、例示および説明の目的で提示されたものである。本発明を網羅すること、または開示された精密な形に限定することは意図されていない。上記の教示に鑑みて、多くの修正および変更が可能である。本発明の範囲は、この詳細な説明によってではなく、添付の特許請求の範囲によって限定されることが意図される。前述の仕様、例、およびデータは、本発明の構成の製造および使用についての完全な説明を提供する。本発明の多くの実施形態は、本発明の趣旨および範囲を逸脱することなく作成可能であるため、本発明は以下に添付の特許請求の範囲内にある。

10

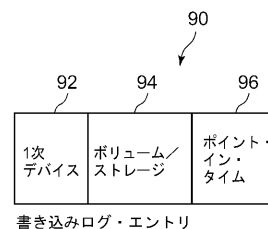
【図 1】



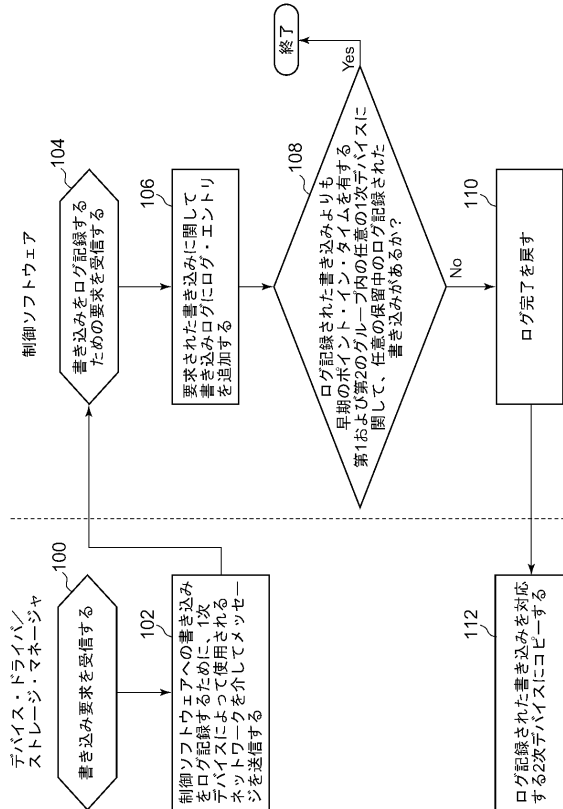
【図 2】



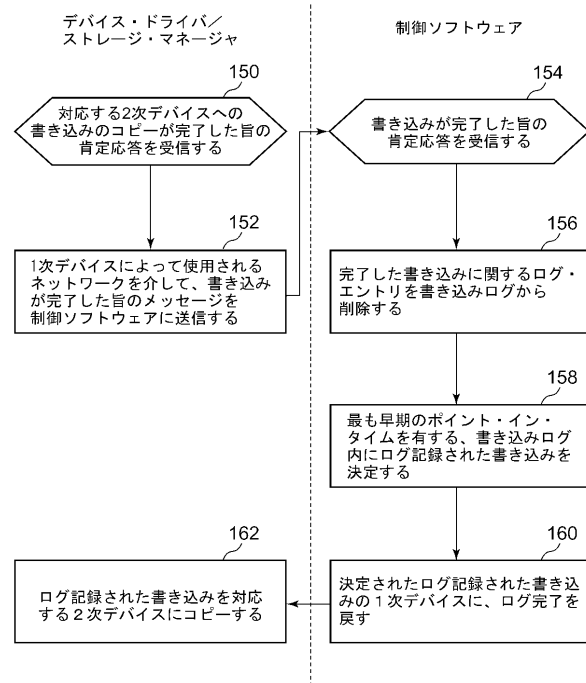
【図 3】



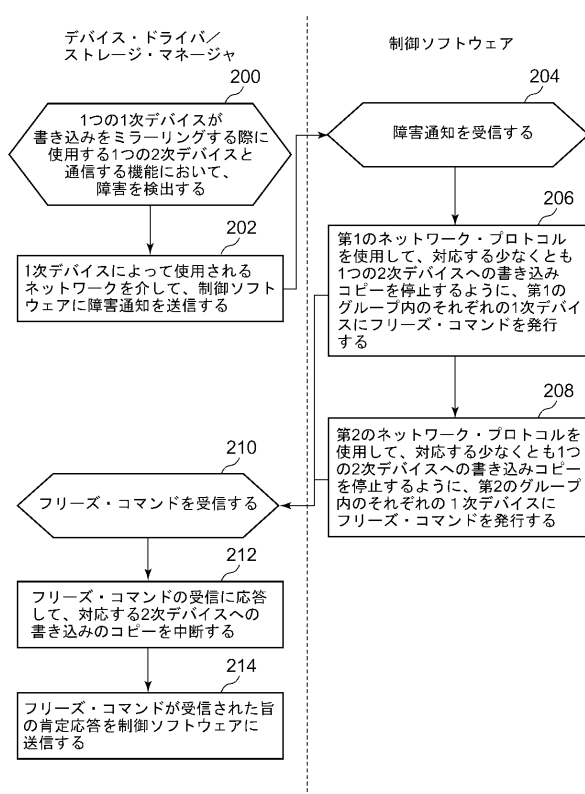
【図 4】



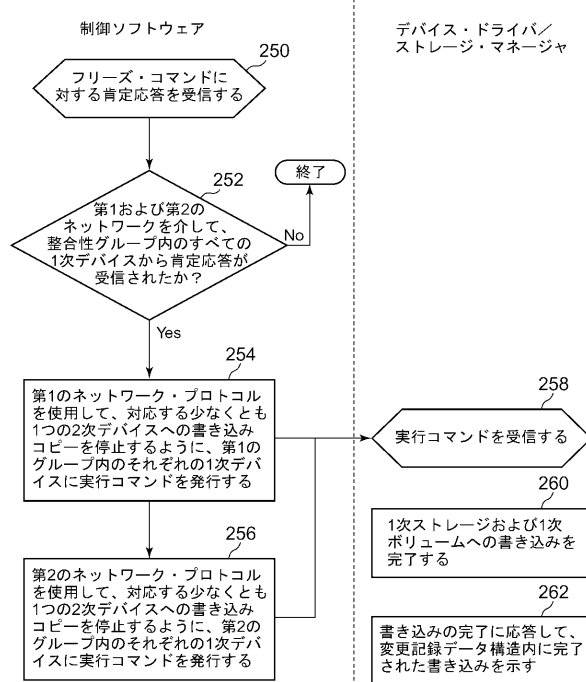
【図 5】



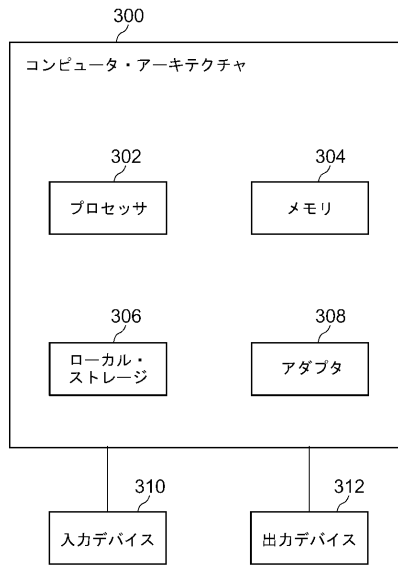
【図 6】



【図 7】



【図 8】





---

フロントページの続き

- (72)発明者 エデル、ステファン、フランス  
アメリカ合衆国 9 6 0 3 7 カリフォルニア州グリーンビュー サウス・キッター・ループ 5 7 0  
4
- (72)発明者 ボイド、ケネス、ワイネ  
アメリカ合衆国 8 5 7 4 8 アリゾナ州ツーソン ノース・カミノ・コードン 8 5 5
- (72)発明者 デイ、サード、ケネス、フェアクロ  
アメリカ合衆国 8 5 7 4 8 アリゾナ州ツーソン ノース・レジュー・ジェイウェイ 7 3 0
- (72)発明者 マクブライド、グレゴリー、エドワード  
アメリカ合衆国 8 5 6 4 1 アリゾナ州ヴェイル サウス・リンコン・ヴァレー・ランチ・ロード  
7 6 5 1

審査官 横山 佳弘

- (56)参考文献 特開 2 0 0 5 - 3 3 2 3 5 4 ( J P , A )  
特開 2 0 0 7 - 1 6 4 7 6 9 ( J P , A )  
特開 2 0 0 6 - 3 1 8 4 6 7 ( J P , A )  
特開 2 0 0 7 - 1 4 1 2 1 6 ( J P , A )  
特開 2 0 0 5 - 1 9 6 6 1 8 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)

G06F 3/06

G06F 12/00

G06F 13/10