

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5159797号
(P5159797)

(45) 発行日 平成25年3月13日 (2013. 3. 13)

(24) 登録日 平成24年12月21日 (2012. 12. 21)

(51) Int. Cl.		F I			
G06F 3/06	(2006.01)	G06F 3/06	302A		
G06F 12/00	(2006.01)	G06F 12/00	531D		
G06F 12/08	(2006.01)	G06F 12/08	541C		
		G06F 12/08	557		

請求項の数 14 (全 16 頁)

(21) 出願番号	特願2009-549389 (P2009-549389)	(73) 特許権者	390009531
(86) (22) 出願日	平成20年1月30日 (2008. 1. 30)		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2010-519607 (P2010-519607A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公表日	平成22年6月3日 (2010. 6. 3)		アメリカ合衆国10504 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/EP2008/051143	(74) 代理人	100108501
(87) 国際公開番号	W02008/101776		弁理士 上野 剛史
(87) 国際公開日	平成20年8月28日 (2008. 8. 28)	(74) 代理人	100112690
審査請求日	平成22年9月30日 (2010. 9. 30)		弁理士 太佐 種一
(31) 優先権主張番号	11/676, 634	(74) 代理人	100091568
(32) 優先日	平成19年2月20日 (2007. 2. 20)		弁理士 市位 嘉宏
(33) 優先権主張国	米国 (US)		
早期審査対象出願			

最終頁に続く

(54) 【発明の名称】 フェイルオーバー後のキャッシュ・データの保存

(57) 【特許請求の範囲】

【請求項1】

データ・ストレージ・サブシステムの一対のクラスタのうちの1つのクラスタであって、前記データ・ストレージ・サブシステムは、少なくとも1つのホスト・アダプタについてのデータを格納するように構成され、ディスク・ストレージは、データを格納するように構成され、前記一対のクラスタは、故障したクラスタと対を成すローカル・クラスタを含み、前記ローカル・クラスタは、

不揮発性データ・ストレージと、
キャッシュ・データ・ストレージと、
サブシステム制御部と

を含み、前記サブシステム制御部は、前記不揮発性データ・ストレージ内に論理サブシステムの書き込みデータを格納し、前記キャッシュ・データ・ストレージ内に別の論理サブシステムの書き込みデータを格納するように構成され、かつ、前記故障したクラスタから前記ローカル・クラスタへのフェイルオーバーにตอบสนองして、前記キャッシュ・データ・ストレージの書き込みデータのトラックの識別情報を、前記ローカル・クラスタのトラックの識別情報に変換することにより、前記書き込みデータを変換済み書き込みデータへと変換し、前記変換済み書き込みデータに対して他の書き込みデータに優先する前記ディスク・ストレージへのデータのデステージの優先順位を与えるように構成され、

前記サブシステム制御部はさらに、前記変換済み書き込みデータの前記ディスク・ストレージへのデステージの失敗にตอบสนองして、不揮発性ストレージのトラックを割り当て、不

揮発性ストレージによる格納のために、前記変換済み書き込みデータを前記不揮発性ストレージに提供し、前記不揮発性ストレージの前記変換済み書き込みデータのトラックの識別情報を、書き込みデータのトラックの識別情報に変換することにより、前記変換済み書き込みデータを、前記不揮発性ストレージ内に格納され、かつ前記キャッシュ・データ・ストレージ内に格納されるローカルDFWデータに再変換するように構成される、クラスタ。

【請求項2】

前記サブシステム制御部は、不揮発性ストレージ制御部及びキャッシュ制御部として具体化され、前記キャッシュ制御部は、前記変換済み書き込みデータが前記ディスク・ストレージへのデステージのための前記優先順位に従って処理されることになるように、前記フェイルオーバにより変換済み書き込みデータについての新規リストを作成するように構成される、請求項1に記載のクラスタ。

10

【請求項3】

前記サブシステム制御部はさらに、前記不揮発性ストレージ内に、前記キャッシュ・データ・ストレージ内に格納された前記1つの論理サブシステムの前記書き込みデータのトラックの識別情報を格納するように構成され、かつ前記故障したクラスタから前記ローカル・クラスタへの前記フェイルオーバに応答して、前記故障したクラスタの前記書き込みデータの前記トラックの識別情報を前記ローカル・クラスタのトラックの識別情報に変換し、これにより前記キャッシュ・データ・ストレージの書き込みデータを変換済み書き込みデータに変換し、変更されたトラックの識別情報を前記新規リストに追加するように構成される、請求項2に記載のクラスタ。

20

【請求項4】

前記サブシステム制御部は、前記提供された変換済み書き込みデータを移動するために、メッセージを前記不揮発性ストレージに送信するように構成される、請求項1に記載のクラスタ。

【請求項5】

前記サブシステム制御部は、前記不揮発性ストレージ内に、前記キャッシュ・データ・ストレージ内に格納された書き込みデータのトラックの識別情報を格納するように構成され、さらに、前記メッセージに対する不揮発性ストレージからの確認信号に応答して、移動された前記変換済み書き込みデータについて、前記不揮発性ストレージから前記トラックの識別情報を削除するように構成される、請求項4に記載のクラスタ。

30

【請求項6】

前記サブシステム制御部は、前記変換済み書き込みデータから再変換された前記書き込みデータを、前記ディスク・ストレージにデステージされるべきDFWリストに追加する、請求項5に記載のクラスタ。

【請求項7】

少なくとも1つのホスト・アダプタについてのデータを格納するように構成されたデータ・ストレージ・サブシステムであって、

データを格納するように構成されたディスク・ストレージと、

各々のクラスタが請求項1から請求項6のいずれか1項に記載の種類である一対のクラスタと

40

を含む、データ・ストレージ・サブシステム。

【請求項8】

データ・ストレージ・サブシステムを含む一対のクラスタのうち故障したクラスタと対を成すローカル・クラスタを動作させるための方法であって、前記データ・ストレージ・サブシステムは、少なくとも1つのホスト・アダプタについてのデータを格納するように構成され、かつデータを格納するように構成されたディスク・ストレージを含み、前記クラスタは、不揮発性データ・ストレージと、キャッシュ・データ・ストレージと、サブシステム制御部とを含み、

前記不揮発性データ・ストレージ内に論理サブシステムの書き込みデータを格納し、前

50

記キャッシュ・データ・ストレージ内に別の論理サブシステムの書き込みデータを格納するステップと、

前記故障したクラスタから前記ローカル・クラスタへのフェイルオーバにตอบสนองして、前記キャッシュ・データ・ストレージの書き込みデータを変換済み書き込みデータに変換し、前記変換済み書き込みデータに対して、他の書き込みデータに優先する前記ディスク・ストレージへのデータのデステージの優先順位を与えるステップと、

前記変換済み書き込みデータの前記ディスク・ストレージへのデステージの失敗にตอบสนองして、不揮発性ストレージのトラックを割り当て、ホスト・アダプタをエミュレートして、前記変換済み書き込みデータを前記不揮発性ストレージによる格納のために前記不揮発性ストレージに提供し、前記不揮発性ストレージの前記変換済み書き込みデータのトラックの識別情報を、書き込みデータのトラックの識別情報に変換することにより、前記不揮発性ストレージの前記変換済み書き込みデータを、前記不揮発性ストレージ内に格納され、かつ前記キャッシュ・データ・ストレージ内に格納される書き込みデータに再変換するステップと

を含む方法。

【請求項 9】

前記変換済み書き込みデータが前記ディスク・ストレージへのデステージのための優先順位に従って処理されることになるように、前記フェイルオーバにより変換済み書き込みデータについての新規リストを作成するステップを含む、請求項 8 に記載の方法。

【請求項 10】

前記不揮発性ストレージ内に、前記キャッシュ・データ・ストレージ内に格納された前記別の論理サブシステムの書き込みデータのトラックの識別情報を格納するステップと、

前記故障したクラスタからローカル・クラスタへの前記フェイルオーバにตอบสนองして、故障したクラスタの前記書き込みデータのトラックの識別情報を前記ローカル・クラスタのトラックの識別情報に変換し、これにより前記キャッシュ・データ・ストレージの書き込みデータを変換済み書き込みデータに変換するステップと、

変更されたトラックの識別情報を前記新規リストに追加するステップと、を含む、請求項 9 に記載の方法。

【請求項 11】

前記提供された変換済み書き込みデータを移動するために、メッセージを前記不揮発性ストレージに送信するステップをさらに含む、請求項 8 に記載の方法。

【請求項 12】

前記サブシステム制御部は、前記不揮発性ストレージ内に、前記キャッシュ・データ・ストレージ内に格納された書き込みデータのトラックの識別情報を格納するように構成され、

前記メッセージに対する前記不揮発性ストレージからの確認信号にตอบสนองして、移動された前記変換済み書き込みデータについて、前記不揮発性ストレージからキャッシュ・データ・ストレージのためのトラックの識別情報を削除するステップをさらに含む、請求項 11 に記載の方法。

【請求項 13】

不揮発性ストレージ・トラックを割り当て、ホスト・アダプタをエミュレートする前記ステップは、

前記不揮発性ストレージを呼び出して、書き込みのための不揮発性ストレージ・セグメントを割り当て、前記ホスト・アダプタをエミュレートして、前記キャッシュ・データ・ストレージから前記割り当てられた不揮発性ストレージ・セグメントに前記変換済み書き込みデータをコピーするステップを含む、請求項 11 に記載の方法。

【請求項 14】

データ・ストレージ・サブシステムを含む一対のクラスタのうち故障したクラスタと対を成すローカル・クラスタを動作させるためのコンピュータ使用可能プログラム・コードを有するコンピュータ・プログラムであって、前記データ・ストレージ・サブシステムは

10

20

30

40

50

、少なくとも1つのホスト・アダプタについてのデータを格納するように構成され、かつデータを格納するように構成されたディスク・ストレージを含み、前記ローカル・クラスタは、不揮発性データ・ストレージと、キャッシュ・データ・ストレージと、サブシステム制御部とを含み、前記サブシステム制御部は、前記不揮発性データ・ストレージ内に論理サブシステムの書き込みデータを格納し、前記キャッシュ・データ・ストレージ内に別の論理サブシステムの書き込みデータを格納し、前記コンピュータ使用可能プログラム・コードは、前記サブシステム制御部上で実行された場合、前記サブシステム制御部に、請求項8～13のいずれか1項に記載のステップを実行するように構成される、コンピュータ・プログラム。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、データ・ストレージ・サブシステムの分野に関し、より具体的には、一対のクラスタのうち的一方からローカル・クラスタへのフェイルオーバーが発生した場合のデータの保存の管理に関する。

【背景技術】

【0002】

データ・ストレージ・サブシステムは、最初にホスト・システムからのDASD高速書き込みデータのようなデータを格納し、次いでDASD又はディスク・ドライブのようなより永続的なデータ・ストレージにデータをデステージするための、種々の形態のデータ・ストレージを含むことができる。一例において、データ・ストレージ・サブシステムは、揮発性であるキャッシュ・データ・ストレージと不揮発性のデータ・ストレージとを各々が有する一対のクラスタを含むことができる。一対のクラスタは、論理サブシステムのデュアル・モードDASD高速書き込みデータの1つのセットが第1のクラスタのキャッシュ・データ・ストレージ及び第2のクラスタの不揮発性データ・ストレージ内に格納され、別の論理サブシステムのデュアル・モードDASD高速書き込みデータが第2のクラスタのキャッシュ・データ・ストレージ及び第1のクラスタの不揮発性データ・ストレージ内に格納されるという点で、データのバックアップを提供する。一例として、デュアル・クラスタ・モードにおいて、偶数の論理サブシステムは、左のクラスタ内のキャッシュと右の不揮発性ストレージとを使用し、奇数の論理サブシステムは、右のクラスタ内のキャッシュと左の不揮発性ストレージとを使用する。

20

30

【0003】

結果として、デュアル・モードDASD高速書き込みデータの全てが、クラスタのうち一方の不揮発性ストア内に格納されることによって、例えば電源の故障又はリブート事象に対して保護される。

【0004】

クラスタのうち一方に障害が発生した場合には、もう一方のクラスタへのフェイルオーバーが実行され、デュアル・モードDASD高速書き込みデータの全ては、1つのセットのデータが不揮発性データ・ストレージ内に格納され、別のセットのデータがキャッシュ・データ・ストレージ内に格納されているので、そのもう一方のクラスタ上で使用可能である。

40

【0005】

従って、一対のクラスタのうち一方からローカル・クラスタへのフェイルオーバーに回答して、ローカル・キャッシュ・ストレージ内に格納された、かつもう一方のクラスタ内の不揮発性ストレージ内に格納されていたことが知られていたデュアル・モードDASD高速書き込みデータが、そのデータの唯一のコピーとなる。キャッシュ・データ・ストレージは揮発性であり、データの唯一のコピーは脆弱なまま残される。

【発明の概要】

【発明が解決しようとする課題】

【0006】

50

一対のクラスタのうちの片方のフェイルオーバが存在し得る場合に、少なくとも1つの
ホスト・アダプタについてのデータを格納するための、データ・ストレージ・サブシステ
ム、データ・ストレージ・サブシステムのクラスタ、コンピュータ・プログラム及び方法
を提供すること。

【課題を解決するための手段】

【0007】

一対のクラスタのうちの片方のフェイルオーバが存在し得る場合に、少なくとも1つの
ホスト・アダプタについてのデータを格納するための、データ・ストレージ・サブシステ
ム、データ・ストレージ・サブシステムのクラスタ、コンピュータ・プログラム及び方法
が提供される。

10

【0008】

データ・ストレージ・サブシステムの実施形態は、データを格納するように構成された
ディスク・ストレージと、一対のクラスタとを含む。クラスタは、ローカル揮発性デー
タ・ストレージと、ローカル・キャッシュ・データ・ストレージと、論理サブシステムの
D A S D高速書き込みデータをローカル揮発性データ・ストレージ内に格納し、別の論
理サブシステムのD A S D高速書き込みデータをローカル・キャッシュ・データストレ
ージ内に格納するように構成されたサブシステム制御部とを含む。

【0009】

一実施形態において、一対のクラスタのうちの一方からローカル・クラスタへのフェイ
ルオーバにตอบสนองして、ローカル・クラスタは、ローカル・キャッシュ・ストレージのデュ
アル・モードD A S D高速書き込みデータを変換高速書き込みデータ (converted fast w
rite data) へと変換し、この変換高速書き込みデータに対して、他の高速書き込みデー
タに優先するディスク・ストレージへのデータのデステージの優先順位を与える。

20

【0010】

更なる実施形態において、クラスタのサブシステム制御部は揮発性ストレージ制御部
及びキャッシュ制御部として具体化され、キャッシュ制御部は、変換高速書き込みデー
タがディスク・ストレージへのデステージのための優先順位に従って処理されることにな
るように、フェイルオーバによる変換高速書き込みデータについての新規リストを作成す
る。

【0011】

更なる実施形態において、クラスタのサブシステム制御部はさらに、ローカル揮発性
ストレージ内に、ローカル・キャッシュ・データ・ストレージ内に格納された別の論理サ
ブシステムのデュアル・モードD A S D高速書き込みデータのトラックIDエントリを格
納するように構成される。一対のクラスタの一方からローカル・クラスタへのフェイ
ルオーバにตอบสนองして、サブシステム制御部は、故障したクラスタのデュアル・モードD A S D
高速書き込みデータのトラックIDエントリをデータのローカル・トラックIDエントリ
に変換し、これによりローカル・キャッシュ・ストレージのデュアル・モードD A S D高
速書き込みデータを変換高速書き込みデータへと変換し、変更されたトラックIDを新規
リストに追加するように構成される。

30

【0012】

別の実施形態において、データ・ストレージ・サブシステムは少なくとも1つのホスト
・アダプタについてのデータを格納するように構成され、データを格納するように構成さ
れたディスク・ストレージと、一対のクラスタとを含む。クラスタは、ローカル揮発性
データ・ストレージと、ローカル・キャッシュ・データストレージと、論理サブシステ
ムのデュアル・モードD A S D高速書き込みデータをローカル揮発性データ・ストレージ内
に格納し、別の論理サブシステムのデュアル・モードD A S D高速書き込みデータをロー
カル・キャッシュ・データ・ストレージ内に格納するように構成されたサブシステム制
御部とを含む。一対のクラスタのうちの一方からローカル・クラスタへのフェイルオーバに
ตอบสนองして、クラスタは、ローカル・キャッシュ・ストレージのデュアル・モードD A S D
高速書き込みデータを変換高速書き込みデータに変換し、ローカル・キャッシュ・ストレ

40

50

ージからディスク・ストレージにデータをデステージしようと試行し、変換高速書き込みデータのディスク・ストレージへのデステージの失敗にตอบสนองして、ローカル不揮発ストレージ・トラックを割り当て、ホスト・アダプタをエミュレートして、ローカル・キャッシュ・データ・ストレージの変換高速書き込みデータをローカル不揮発性ストレージによる格納のためにローカル不揮発性ストレージに提供するように構成され、ローカル・キャッシュ・データ・ストレージの変換高速書き込みデータを、ローカル不揮発性ストレージ内に格納され、かつローカル・キャッシュ・ストレージ内に格納されるローカル・シングル・モードDASD高速書き込みデータに再変換する。

【0013】

更なる実施形態において、各クラスタのサブシステム制御部は、提供された変換高速書き込みデータをコミットするために、ホスト・アダプタ・スタイル・コミット・メッセージをローカル不揮発性ストレージに送り出すように構成される。

10

【0014】

更なる実施形態においては、クラスタのサブシステム制御部は、ローカル不揮発性ストレージ内に、ローカル・キャッシュ・データ・ストレージ内に格納された高速書き込みデータのトラックIDエントリを格納するように構成され、さらに、ホスト・アダプタ・スタイル・コミット・メッセージに対するローカル不揮発性ストアからのコミット確認信号にตอบสนองして、コミットされた変換高速書き込みデータについて、ローカル不揮発性ストアからキャッシュ・トラックIDエントリを削除するように構成される。

【0015】

20

更なる実施形態において、クラスタのサブシステム制御部は、変換高速書き込みデータから再変換されたDASD高速書き込みデータを、ディスク・ストレージにデステージされるべきDFWリストに追加する。

【0016】

別の実施形態において、クラスタのサブシステム制御部は、不揮発性ストレージ制御部と、キャッシュ制御部と、逆行ストア(retro-store)制御部として具体化される。「逆行ストア制御部」は、ホスト・アダプタをエミュレートする制御部又は制御コードである。キャッシュ制御部は、不揮発性ストレージ制御部を呼び出して、書き込みのための不揮発性ストレージ・セグメントを割り当て、逆行ストア制御部は、ホスト・アダプタをエミュレートして、変換高速書き込みデータを不揮発性ストレージ・セグメントにコピーする。

30

【0017】

ここで、添付の図面を参照して、例示のみの目的で本発明の実施形態を説明する。

【図面の簡単な説明】

【0018】

【図1】本発明の実施形態を実施することができるデータ・ストレージ・サブシステムを示すブロック図である。

【図2】図1のデータ・ストレージ・サブシステムによって格納されるデータ・タイプの、従来技術の一時的格納を説明する図であり、これは逐次高速書き込みデータと呼ぶことができる。

40

【図3】図1のデータ・ストレージ・サブシステムによって格納される別のデータ・タイプの、従来技術の一時的格納を説明する図であり、これはDASD高速書き込みデータと呼ぶことができる。

【図4】図3のDASD高速書き込みデータを受け取り、そのデータを一時的に格納するための従来技術のプロセスを示す図である。

【図5】本発明によるフェイルオーバー処理を示すフローチャートである。

【図6】図5のフェイルオーバー・プロセスを説明する図である。

【図7】図1のデータ・ストレージ・サブシステムのシングル・クラスタ・モードでDASD高速書き込みデータを受け取るための従来技術のプロセスを説明する図である。

【図8】ディスク・ストレージへのデータのデステージ不能のときの、本発明によるフェ

50

イルオーバ処理を示すフローチャートである。

【発明を実施するための形態】

【0019】

以下の説明において、本発明は図面を参照して好ましい実施形態で説明され、ここで同様の数字は同一又は類似の要素を表す。本発明は、本発明の目的を達成するための最良の形態によって説明されるが、当業者であれば、これらの教示を鑑みて、本発明の精神又は範囲から逸脱することなくバリエーションを達成することができることを認識するであろう。

【0020】

図1を参照すると、データ・ストレージ・サブシステム100は、クラスタ110と、
もう一つのクラスタ120とを含む。クラスタ110は、少なくともサブシステム制御部
132とローカル不揮発性データ・ストレージ134とローカル・キャッシュ・データ・
ストレージ136とを組み入れた、複合部130を含む。同様に、クラスタ120は、少
なくともサブシステム制御部142とローカル不揮発性データ・ストレージ144とロー
カル・キャッシュ・データ・ストレージ146とを組み入れた、複合部140を含む。各
クラスタにおいて、サブシステム制御部は、複合部の残りの部分から完全に分離して
もよく、或いはローカル不揮発性データ・ストレージ及び/又はローカル・キャッシュ・
データストレージに部分的に組み入れられていてもよい。サブシステム制御部132、1
42は、論理及び/又は1つ又は複数のマイクロプロセッサを含み、情報及びマイクロ
プロセッサを動作させるためのプログラム情報を格納するためのメモリを備える。本明細書
において、「プロセッサ」又は「制御部」は、任意の適切な論理、プログラム可能論理、
マイクロプロセッサ、及びプログラム命令に応答するための連結されたメモリ又は内部メ
モリを含むことができ、連結されたメモリ又は内部メモリは、固定メモリ若しくは書き換
え可能メモリ又はデータ・ストレージ装置を含むことができる。プログラム情報は、ホス
トから又はデータ・ストレージ・ドライブ若しくはディスク・アレイを介して、又はフロ
ッピー若しくは光ディスクからの入力によって、又はカートリッジからの読み出しによ
って、又はウェブ・ユーザ・インターフェース若しくは他のネットワーク接続によっ
て、又は他のいずれかの適切な手段によって、サブシステム制御部又はメモリに供給す
ることができる。従って、プログラム情報は、クラスタ110を動作させるため及び/又はクラ
スタ120を動作させるための有形に具現化されたコンピュータ使用可能なプログラム・コ
ードをその中に有するコンピュータ使用可能媒体を含む、1つ若しくは複数のプログラム
、又は同様のタイプのシステム若しくはデバイスを含むことができる。

【0021】

不揮発性データ・ストレージ134、144は、当業者に公知のように、電力が失われ
た場合でもデータを保護するバッテリー・バックアップを有するメモリ・システム、フラッ
シュPROM、ディスク・ドライブ、又は他の適切な不揮発性メモリを含むことができ
る。キャッシュ・データ・ストレージ136、146は、当業者に公知のように、任意の適
切なメモリ・システムを含むことができ、揮発性であってもよく、電力が取り去られた後
でデータを失う可能性がある。

【0022】

アダプタ・インターフェース(AI)138、148は、キャッシュ・データ・ストレ
ージ136、146の一部を含むことができ、及び/又はサブシステム制御部132、1
42の一部を含むことができ、かつキャッシュ・データ・ストレージ136、146のと
ころに常駐することもできるし、又は別個に常駐することもでき、若しくは複合体130
、140のその他の要素と共に常駐することもできる。アダプタ・インターフェースは、
当業者に公知のように、特定のクラスタについて、ローカル不揮発性データ・ストレ
ージ及びキャッシュ・データ・ストレージに関するデータ転送のハンドリング面のための論理
を提供する。

【0023】

複数のホスト・アダプタ150-157は、全て当業者に公知のような、1つ又は複数

10

20

30

40

50

のファイバチャネル (Fibre Channel) ポート、1つ又は複数の F I C O N ポート、1つ又は複数の E S C O N ポート、1つ又は複数の S C S I ポート、又は他の適切なポートを含むことができる。各ホスト・アダプタは、各クラスタがどのアダプタからの I / O でも処理することができるように、ホスト・システムと通信し、かつクラスタ 1 1 0 及びクラスタ 1 2 0 の両方と通信するように構成される。

【 0 0 2 4 】

複数のデバイス・アダプタ 1 6 0 - 1 6 7 は、ディスク・アレイ 1 7 0 - 1 7 3 のようなディスク・ドライブ又はディスク・ドライブ・システムと通信するための、通信リンクを含むことができる。或いは、磁気テープ・ドライブで1つ又は複数のディスク・アレイを置き換えることができる。ディスク・アレイは、R A I D (Redundant Array of Independent Disks) プロトコルを利用することもできるし、又は J B O D (Just a Bunch of Disks) アレイを含むこともできる。通信リンクは、R S - 2 3 2 又は R S - 4 2 2 のようなシリアル相互接続、イーサネット接続、S C S I 相互接続、E S C O N 相互接続、F I C O N 相互接続、ローカル・エリア・ネットワーク (L A N)、私設広域ネットワーク (W A N)、公共広域ネットワーク、ストレージ・エリア・ネットワーク (S A N)、伝送制御プロトコル / インターネット・プロトコル (T C P / I P)、インターネット、及びこれらの組み合わせを含むことができる。

10

【 0 0 2 5 】

データ・ストレージ・サブシステム 1 0 0 の例は、I B M (登録商標) エンタープライズ・ストレージ・サーバ (Enterprise Storage Server) モデル D S / 8 0 0 0、又はその他の相当するシステムを含む。

20

【 0 0 2 6 】

上述のように、データ・ストレージ・サブシステムは、ホスト・システムからのデータを格納するための種々の形態のデータ・ストレージを含むことができる。

【 0 0 2 7 】

図 1 及び図 2 を参照すると、1つのタイプのホスト・データは、データ・ストレージ・サブシステム 1 0 0 によって格納される一方で、システム内のどこか他のポイントにオリジナルのコピーが存在する、逐次高速書き込み (S F W) データである。例えば、データが一次サイトに格納され、これが二次サイトとしてのデータ・ストレージ・サブシステムのディスク・アレイ 1 7 0 - 1 7 3 にコピーされる場合には、データは、一次制御部から到着すべき二次サイトへの「 P P R C 確立 (Establish) 」 (ピアツーピア遠隔コピー (Peer-to-Peer RemoteCopy)) の間に、ホスト・アダプタ 1 5 8 を介して S F W データとして送られる。ホスト・アダプタ 1 5 8 は、ホスト・アダプタ 1 5 0 - 1 5 7 のうちのいずれか又は1つより多くについての代理として図示されている。このデータは、キャッシュ・データ・ストレージ 1 3 6、1 4 6 にのみ送られ、不揮発性ストレージ 1 3 4、1 4 4 には送られない。トラック I D (識別) エントリが、トラック毎に不揮発性ストレージ 1 3 4、1 4 4 内に置かれる。このエントリは、キャッシュ情報の損失を引き起こすリポート動作が生じた場合に必要とされる。データはその後、一次制御部から再度受け取ったときにキャッシュ・ストレージ内で正確に復旧されることができる。この例において、ステップ 1 8 0 で、S F W データがホスト・アダプタ 1 5 8 によってクラスタ 1 2 0 のキャッシュ・データ・ストレージ 1 4 6 内に格納され、ステップ 1 8 1 で、キャッシュがトラック I D エントリをクラスタ 1 1 0 の不揮発性ストレージ 1 3 4 内に格納する。一旦、全てのデータが一次から二次にコピーされると、二次への更なる書き込みは全て、D A S D 高速書き込み (D F W) データとして到着する。1つのクラスタからローカル・クラスタへのフェイルオーバが生じた場合には、ローカル・クラスタの不揮発性ストレージ内のトラック I D エントリを用いて、故障したクラスタの S F W データに再アクセスする。

30

40

【 0 0 2 8 】

図 1 及び図 3 を参照すると、別のタイプのホスト・データは D A S D 高速書き込み (D F W) データであり、これは最初にホスト・システムから 1 5 8 を介してクラスタ 1 1 0 及び 1 2 0 によって格納され、その後、D A S D 又はディスク・ドライブ 1 7 0 - 1 7 3

50

のような、より永続的なデータ・ストレージにデステージされる。一対のクラスタ110及びクラスタ120は、1つの論理サブシステムのデュアル・モードDASD高速書き込みデータの1つのセットが第1のクラスタのキャッシュ・データ・ストレージ136及び第2のクラスタの不揮発性データ・ストレージ144内に格納され、別の論理サブシステムのデュアル・モードDASD高速書き込みデータが第2のクラスタのキャッシュ・データ・ストレージ146及び第1のクラスタの不揮発性データ・ストレージ134内に格納されるという点で、データのバックアップを提供する。一例として、デュアル・クラスタ・モードにおいて、偶数の論理サブシステムは、左のクラスタのキャッシュと右の不揮発性ストレージとを使用し、奇数の論理サブシステムは、右のクラスタのキャッシュと左の不揮発性ストレージとを使用する。それに加えて、キャッシュ・トラックIDエントリが、他方のクラスタの不揮発性ストレージ内に格納される。この例において、DFWデータは、ステップ185で、ホスト・アダプタ158によってクラスタ120のキャッシュ・データ・ストレージ146内に格納され、ステップ186で、クラスタ110の不揮発性ストレージ134内に格納される。それに加えて、クラスタ120のキャッシュ146は、クラスタ110の不揮発性ストレージ134内にキャッシュ・トラックIDエントリを格納する。

10

【0029】

結果として、デュアル・モードDASD高速書き込みデータの全てが、クラスタのうちの1つの中の不揮発性ストア内に格納されることによって、例えば、電源異常又はリポート事象から保護される。

20

【0030】

図4は、その後でデステージするためにデュアル・モードDASD高速書き込みデータを格納する従来技術の詳細なプロセスの例を示す。

【0031】

図1、図3及び図4を参照すると、ステップ201において、ホスト・アダプタ158は、トラックに対する書き込み要求を取得し、ステップ202において、クラスタ「B」のキャッシュ146のアダプタ・インターフェース148に、キャッシュ及び不揮発性ストレージ(NVS)セグメント並びに不揮発性ストレージ・バッファを割り当てるためのトラックIDと共に、メールを送信する。ステップ203において、アダプタ・インターフェース148は、キャッシュ146を呼び出して、キャッシュ/NVSセグメントを割り当て、キャッシュ制御ブロックを作成する。アダプタ・インターフェースはまた、NVSトラック・バッファも割り当てる。ステップ204において、アダプタ・インターフェースは、トラック制御ブロックを構築し、それをクラスタ「A」のNVS134に送って、使用すべきセグメントを指定する。ステップ205において、アダプタ・インターフェースは、ホスト・アダプタにNVSトラック・バッファ番号と共にメールを送信して、書き込みを開始させる。ステップ206において、ホスト・アダプタ158は、DMA機能(直接メモリ・アクセス)を使用して、データをキャッシュ・セグメント及びNVSトラック・バッファに送り、ステップ207(同じ矢印)において、データをコミットするために、キャッシュにメールを送信し、データをコミットするのに使用されたトラック・バッファ番号と共に、NVSにメールを送信する。ステップ208において、ホスト・アダプタは、書き込みの完了を示すデバイス・エンドをホスト・システムに与える。書き込みの完了は、クラスタへの電力が失われた後でさえもNVS134がデータをコミットするという事実によってサポートされる。ステップ209において、NVS134は、メール207を確認し、次いでトラックについてのNVS制御ブロックを構築し、トラック・バッファからNVSセグメントにデータを移動させることによって、データをコミットし、ステップ210において、NVSは、アダプタ・インターフェース148にコミット完了のメールを送信する。ステップ211において、アダプタ・インターフェースは、ホスト・アダプタとNVSの両方からの「完了」メールを確認し、ステップ212において、キャッシュ146を呼び出して、キャッシュとNVSの両方に書き込まれたセグメントによってキャッシュ制御ブロックを更新し、NVSトラック・バッファを解放する。ステップ

30

40

50

213において、アダプタ・インターフェースは、書き込み完了のメッセージをホスト・アダプタ158に送る。このようにして、ホスト・アダプタは、ステップ208のデバイス・エンドが完全にサポートされ、デュアル・モードDASD高速書き込みデータがクラスタ「A」のNVS134とクラスタ「B」のキャッシュ146の両方に書き込まれ、格納されていることを知る。他の詳細なシーケンスを用いて、その後でデステージするためのデュアル・モードDASD高速書き込みデータの格納を達成することができる。クラスタ「B」のキャッシュ146はまた、上述のように、データ・ストレージ・サブシステムのデュアル・モード逐次高速書き込み(SFW)データの半分も含む。

【0032】

一例において、奇数の論理サブシステムからのデュアル・モードDASD高速書き込みデータは、クラスタ「A」のNVS134内に格納され、(上述の)奇数の論理サブシステムからのデュアル・モードDASD高速書き込みデータ及びデュアル・モード逐次高速書き込みデータは、両方ともクラスタ「B」のキャッシュ146内に格納される。同様に、偶数の論理サブシステムからのデュアル・モードDASD高速書き込みデータは、クラスタ「B」のNVS144内に格納され、偶数の論理サブシステムからのデュアル・モードDASD高速書き込みデータ及びデュアル・モード逐次高速書き込みデータは、両方ともクラスタ「A」のキャッシュ136内に格納される。

10

【0033】

キャッシュ136、146は、典型的には、不揮発性ストレージ134、144よりも格納されるデータ量当たりの費用が低いので、従って、逐次高速書き込みデータとDASD高速書き込みデータの両方を処理するために、より大きな容量で提供される。

20

【0034】

一方のクラスタから他方のクラスタへのフェイルオーバーが生じた場合、デュアル・モードDASD高速書き込みデータの全てが他方のクラスタ上で使用可能であり、デュアル・モードDASD高速書き込みデータの1つのセットは不揮発性データ・ストレージ内に格納され、データのもう1つのセットは、逐次高速書き込みデータのような他の高速書き込みデータと共にキャッシュ・データ・ストレージ内に格納される。

【0035】

フェイルオーバーの結果として、ローカル・キャッシュ・ストレージ内に格納された、かつ他方のクラスタ内の不揮発性ストレージ内に格納されていたことが知られていたデュアル・モードDASD高速書き込みデータが、そのデータの唯一のコピーになる。キャッシュ・データ・ストレージは揮発性であり、DASD高速書き込みデータの唯一のコピー及び他の高速書き込みデータは脆弱なまま残される。従って、典型的には、キャッシュ・データ・ストレージ内にあるデータは、データを保護するために、ディスク・ストレージ170-173のような、より永続的なストレージにデステージされる。

30

【0036】

本発明によれば、図1、図5及び図6を参照すると、ステップ240における一对のクラスタのうち的一方からローカル・クラスタへのフェイルオーバーにตอบสนองして、一実施形態において、ローカル・クラスタは、ローカル・キャッシュ・ストレージのデュアル・モードDASD高速書き込みデータを変換高速書き込みデータへと変換し、この変換高速書き込みデータに、他の高速書き込みデータに優先するディスク・ストレージへのデータのデステージの優先順位を与える。この例においては、クラスタ110が故障したものと仮定し、フェイルオーバーはクラスタ120へのフェイルオーバーである。図6は、不揮発性データ・ストレージ144並びにキャッシュ146及びアダプタ・インターフェース148のみを図示するものであり、サブシステム制御部142又はクラスタの他の局面は図示せず、ディスク・ストレージ170-173を、ディスク・ストレージ170-173の代理として図示されるディスク・ストレージ174として特徴付ける。

40

【0037】

一実施形態において、ステップ243において、サブシステム制御部142は、不揮発性ストレージ144にアクセスして、故障したクラスタのデュアル・モードDASD高速

50

書き込みトラックIDエントリのリストを提供する。このリストは、上述のようにNV S 1 4 4によって提供される高速書き込みデータの完全なリストから切り離され、キャッシュ1 4 6に提供することができる。ステップ2 4 5において、ストレージ制御部は、NV S 1 4 4を動作させて、例えばキャッシュ1 4 6に、トラック制御ブロックを提供する。ステップ2 4 7において、サブシステム制御部は、トラック制御ブロックを処理して、故障したクラスタのデュアル・モードD A S D高速書き込みデータのトラックIDエントリをデータのローカル・トラックIDエントリに変換し、これによりローカル・キャッシュ・ストレージのデュアル・モードD A S D高速書き込みデータを変換高速書き込みデータへと変換する。ステップ2 4 9において、サブシステム制御部は、NV Sについて、LR U（最長時間未使用）リストのような新規リストを作成し、変換高速書き込みデータのIDエントリをその新規リストに追加する。或いは、新規リストは、F I F O（先入れ先出し）リストを含むことができる。一実施形態において、トラック制御ブロック及び変換トラックIDエントリを処理するためのサブシステム制御部コードは、キャッシュ1 4 6のためのキャッシュ制御モジュールと共に常駐する。代替的な実施形態においては、サブシステム制御部コードの少なくとも一部は、プロセッサ1 4 2及びキャッシュ制御部1 4 6から分離しており、アダプタ・インターフェース1 4 8と共に常駐する。さらに、サブシステム制御部コードの少なくとも一部は、不揮発性ストレージ1 4 4のための不揮発性ストレージ制御モジュールと共に常駐することができる。従って、一実施形態において、アダプタ・インターフェース1 4 8がキャッシュ制御部を呼び出し、フェイルオーバによる変換高速書き込みデータについての新規リストを作成する。ステップ2 5 0において、サブシステム制御部は、例えば、新規リストのLR U又はF I F Oの基準に基づいて、変換高速書き込みデータ・トラックの新規リストに対してディスク・ストレージへのデステージのための優先順位を与える。変換高速書き込みデータ・トラックが首尾良くデステージされた場合、NV S内のトラックIDエントリが削除され、トラックは、キャッシュ内の未修飾トラックとなる変更される。

【0038】

従って、キャッシュ・データ・ストレージ内の変換高速書き込みデータが、ディスク・ストレージのような、より永続的なストレージにデステージされる。

【0039】

それに加えて図7を参照すると、このデータ・ストレージ・システムは、その後、残りのクラスタ1 2 0のみを用いて、従来技術のシングル・モード・データ・ストレージ・サブシステムとして動作を続けることができ、そこでは、ホスト・アダプタ1 5 8は、デステージされるべきD F Wデータを同じクラスタのNV S 1 4 4とキャッシュ1 4 6の両方に提供し、データ・ストレージの1つのモードが他方に対するバックアップとして機能するので、ある程度の安全性を有する。

【0040】

通常、デステージ・プロセスは迅速に実行されるが、データ・ストレージ・システムは、デステージの成功を阻害するドライブ若しくはランク、又はデバイス・アダプタの問題に遭遇する可能性がある。その結果、ローカル・キャッシュ内にのみ格納されているデュアル・モードD A S D高速書き込みデータの一部は脆弱なまま残され、データの唯一のコピーとなる。

【0041】

図1、図6及び図8を参照すると、本発明によれば、一对のクラスタのうち的一方からローカル・クラスタへのフェイルオーバ、例えばクラスタ1 2 0へのフェイルオーバに回答して、クラスタは、ローカル・キャッシュ・データ・ストレージからディスク・ストレージにデータをデステージしようと試行するように構成される。ステップ3 0 0の変換高速書き込みデータのディスク・ストレージへのデステージの失敗に回答して、クラスタは、ローカル不揮発性ストレージ・トラックを割り当て、ホスト・アダプタをエミュレートして、ローカル・キャッシュ・データ・ストレージ1 4 6の変換高速書き込みデータをローカル不揮発性ストレージによる格納のためにローカル不揮発性ストレージ1 4 4に提供

10

20

30

40

50

して、ローカル・キャッシュ・データストレージの変換高速書き込みデータを、ローカル不揮発性ストレージ内に格納され、かつローカル・キャッシュ・ストレージ内に格納されるローカル・シングル・モードDASD高速書き込みデータに再変換する。

【0042】

一実施形態において、ステップ303で、ストレージ制御部は、ホスト・アダプタ158のようなホスト・アダプタをエミュレートし、また、アダプタ・インターフェース148のようなアダプタ・インターフェースをエミュレートすることもできる。コードは、「逆行ストア」制御部と呼ぶことができる。ホスト・アダプタをエミュレートするストレージ制御部は、DFW書き込み動作のためのトラックにアクセスするためにキャッシュ又は実際のアダプタ・インターフェース若しくはエミュレートされたアダプタ・インターフェースを呼び出し、キャッシュ等は書き込みのためにローカルNVS144のNVSセグメントを割り当てる。例えば、エミュレートされたアダプタ・インターフェース・コードはキャッシュを呼び出してトラックのためのキャッシュ制御ブロックをロックし、キャッシュは、書き込みのために割り当てられるNVS空間を得るために、NVSに対してインターフェースする。例えばアダプタ・インターフェースをエミュレートするストレージ制御部は、ローカルNVS144内にNVSトラック・バッファを割り当てることもできる。ストレージ制御部は、トラックNVS制御ブロックを構築し、それをローカルNVSに送って、使用すべきセグメントを指示する。例えば、キャッシュは、空間が割り当てられていることを示すエミュレートされたアダプタ・インターフェース・コードへと戻り、このアダプタ・インターフェース・コードがNVSを呼び出して、書き込みを開始するためにNVSトラック・バッファを割り当てる。

10

20

【0043】

ホスト・アダプタをエミュレートするストレージ制御部は、上記の変換高速書き込みデータを、ローカル不揮発性ストレージによる格納のためにローカル不揮発性ストレージ144に、例えば上記のNVSトラック・バッファに、コピーする。

【0044】

一実施形態において、ステップ305で、ホスト・アダプタをエミュレートするストレージ制御部の逆行ストアコードは、ホスト・アダプタ・スタイル・コミット・メッセージをローカルNVSに送り、例えば、トラックについてコピーされたデータをコミットするためにローカル不揮発性ストア144にメールを送信する。ステップ307において、ローカルNVSはメールを確認し、次に、例えばトラックについてのNVS制御ブロックを構築し、コピーされたデータをNVSトラック・バッファから上記で割り当てられたNVSセグメントに移動させることによって、コピーされたデータをコミットする。NVSは次に、コミット完了メール・メッセージを逆行ストアコードに送り返し、コミット確認信号(acknowledge)を提供する。

30

【0045】

ステップ309において、ストレージ制御部の逆行ストアは、コミット完了メール・メッセージに回答して、NVSトラック・バッファを解放し、次いでキャッシュを呼び出して、キャッシュ146のためのトラックIDエントリをローカルNVSから削除する。

【0046】

上述のようなサブシステム制御部は多くの形態をとることができ、ローカル・キャッシュ・データ・ストレージ内に格納された高速書き込みデータのトラックIDエントリをローカル不揮発性ストレージ内に格納するように構成され、さらに、ローカル不揮発性ストアからホスト・アダプタ・スタイル・コミット・メッセージへのコミット確認信号に回答して、コミットされた変換高速書き込みデータについてのキャッシュ・トラックIDエントリをローカル不揮発性ストアから削除するように構成される。

40

【0047】

更なる実施形態においては、ステップ311において、一旦、トラックIDエントリが削除されると、ストレージ制御部の逆行ストアは、トラックを再変換して(上述の)変換高速書き込みデータ・トラックからDASD高速書き込みデータ・トラックに戻すが、こ

50

のときには、シングル・モードDASD高速書き込みデータ・トラックである。例えばアダプタ・インターフェースをエミュレートするストレージ制御部は、キャッシュを呼び出し、トラックを、新規リストから、その機能が復旧したときにディスク・ストレージにデステージされるべきDASD高速書き込み(DFW)トラック・リストに移動させる。

【0048】

当業者であれば、ステップの順序の変更を含めた変更を上述の方法に関して行うことができることを理解するであろう。さらに、当業者であれば、本明細書において例示されたものとは異なる特定のコンポーネント配置を用いることができることを理解するであろう。

【0049】

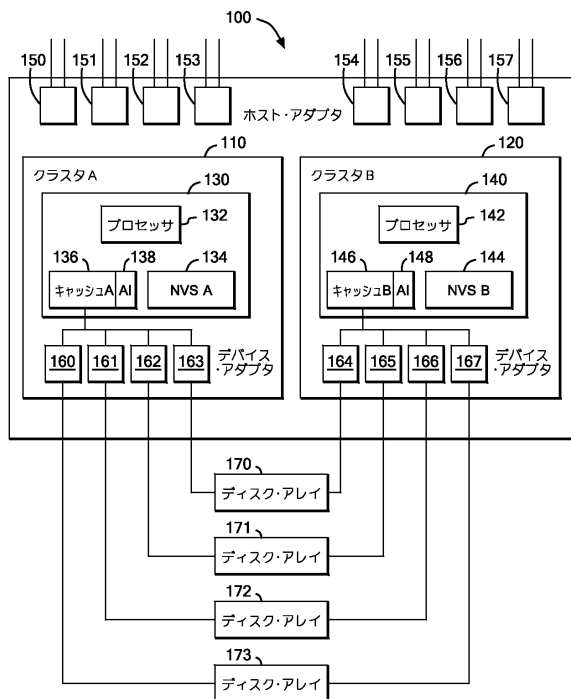
本発明の好ましい実施形態を詳細に例示してきたが、当業者であれば、以下の特許請求の範囲に示される本発明の範囲を逸脱することなく、それらの実施形態への修正及び改変を想起できることが明らかである。

【符号の説明】

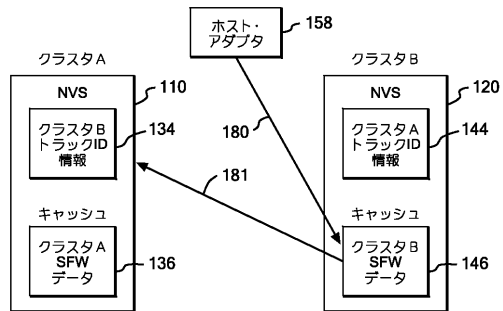
【0050】

- 100：データ・ストレージ・サブシステム
- 110、120：クラスタ
- 130、140：複合部
- 132、142：サブシステム制御部
- 134、144：ローカル不揮発性データ・ストレージ
- 136、146：ローカル・キャッシュ・データ・ストレージ
- 138、148：アダプタ・インターフェース
- 150 - 157、158：ホスト・アダプタ
- 160 - 167：デバイス・アダプタ
- 170 - 173、174：ディスク・ストレージ

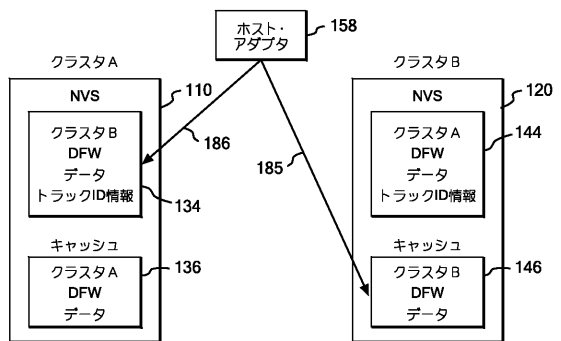
【図1】



【図2】



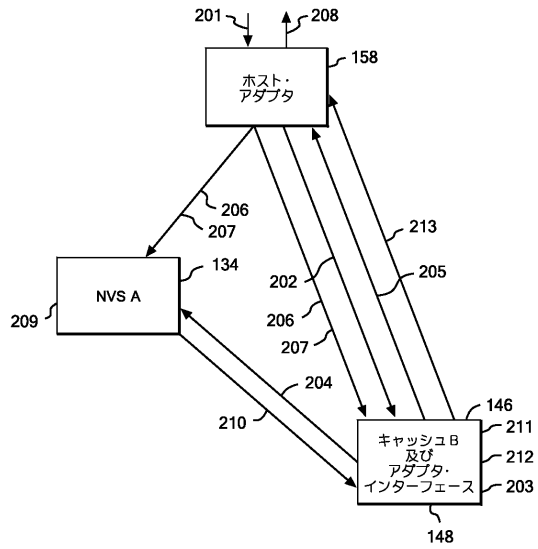
【図3】



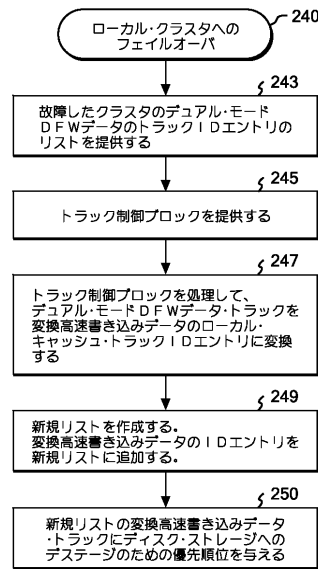
10

20

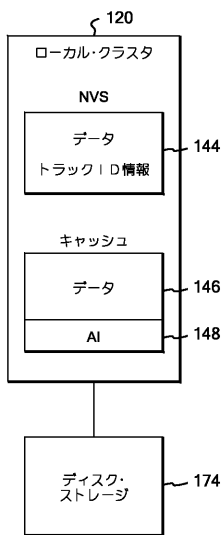
【図4】



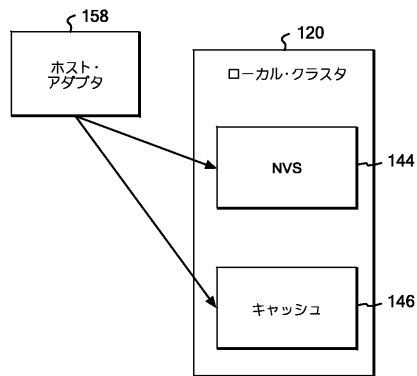
【図5】



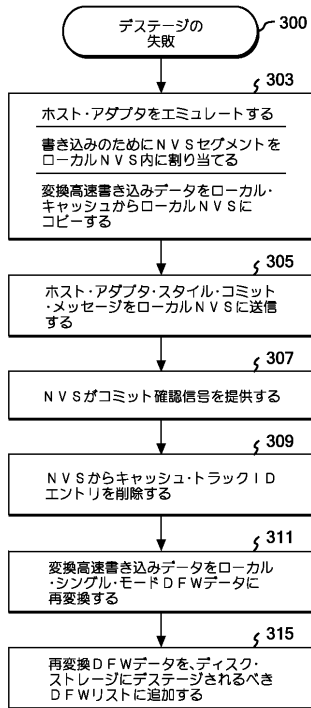
【図6】



【図7】



【図 8】



フロントページの続き

- (72)発明者 アッシュ、ケビン、ジョン
アメリカ合衆国 85746 アリゾナ州 ツーソン サウス・サドル・リッジ・レーン 602
1
- (72)発明者 グプタ、ローケシュ、モハン
アメリカ合衆国 85747 アリゾナ州 ツーソン サウス・リバー・ウィロー・ドライブ 7
452
- (72)発明者 ロウ、スティーブン、ロバート
アメリカ合衆国 85710 アリゾナ州 ツーソン サウス・ヘルモサ・ヒルズ・プレイス 8
31
- (72)発明者 サンチェス、アルフレッド、エミリオ
アメリカ合衆国 85743 アリゾナ州 ツーソン ノース・ミスティ・ブルック・ドライブ
8885
- (72)発明者 トッド、ケネス、ウェイン
アメリカ合衆国 85748 アリゾナ州 ツーソン ノース・イースタン・スロープ・ループ
145

審査官 坂東 博司

- (56)参考文献 特開平06-222988(JP,A)
米国特許第05437022(US,A)
米国特許第05771367(US,A)
特開平08-115169(JP,A)
米国特許出願公開第2004/0059870(US,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06
G06F 12/00
G06F 12/08