

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2006-209775

(P2006-209775A)

(43) 公開日 平成18年8月10日(2006.8.10)

(51) Int. Cl.		F I		テーマコード (参考)
<b>G06F 3/06 (2006.01)</b>		G06F 3/06	304P	5B065
<b>G06F 12/00 (2006.01)</b>		G06F 3/06	540	5B082
		G06F 3/06	305C	
		G06F 12/00	531M	

審査請求 有 請求項の数 10 O L (全 23 頁)

(21) 出願番号	特願2006-17315 (P2006-17315)	(71) 出願人	503003854 ヒューレット・パカード デベロップメント カンパニー エル. ピー. アメリカ合衆国 テキサス州 77070 ヒューストン 20555 ステイト ハイウェイ 249
(22) 出願日	平成18年1月26日 (2006.1.26)	(74) 代理人	110000039 特許業務法人アイ・ピー・エス
(31) 優先権主張番号	11/046350	(72) 発明者	カール・ドーム アメリカ合衆国カリフォルニア州 パロアルト ハノーバー・ストリート 3000 ヒューレット・パカード・カンパニー 内
(32) 優先日	平成17年1月28日 (2005.1.28)	F ターム (参考)	5B065 BA01 CA12 CA30 CC08 EA02 EA12 EA24 EK05 5B082 DA02 DC00 DE05 HA01
(33) 優先権主張国	米国 (US)		

(54) 【発明の名称】 データ追跡を有するストレージ複製システム

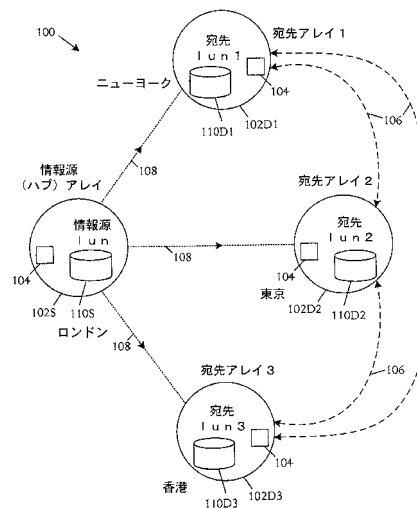
(57) 【要約】

【課題】 データ追跡を備えるストレージ複製システムを提供する。

【解決手段】

本発明にかかるストレージ複製方法は、複数の宛先の中の情報源からデータを複製すること、及び、それら宛先のデータの変更を追跡することを含む。変更の識別情報は、複数の宛先アレイ間で相互に通信される。情報源のフェイルオーバーイベントでは、選択された宛先が新たな情報源として設定され、残りの宛先における複製されたデータを、当該新たな情報源と同期した状態に再形成する。

【選択図】 図1 A



## 【特許請求の範囲】

## 【請求項 1】

冗長データストレージシステム(200)での使用に適合されたストレージユニット(202)であって、

ストレージ(204)と、

前記ストレージ(204)に接続されて、複数の分散サイトストレージユニットと通信するように適合された入出力インターフェース(206)と、

前記ストレージ(204)及び前記入出力インターフェース(206)に接続されたコントローラ(208)と、

前記コントローラ(208)において実行可能なロジックであって、前記ストレージユニット(202)をファンアウト配置における2次サイトとして動作させ、1次サイトストレージユニットから前記ストレージ(204)へデータを複製し、前記ストレージ(204)に書き込まれたデータの変更を追跡し、前記追跡された変更を前記複数の分散サイトストレージユニット間で通信し、前記複数の分散サイトストレージユニットから受信された、追跡された変更を収集するように適合されたロジックと

10

を備えるストレージユニット。

## 【請求項 2】

ストレージユニット(202)であって、

フェイルオーバーするコマンドを受信し、前記フェイルオーバーするコマンドに回答して、前記ストレージユニット(202)を1次サイトとして動作させ、前記複数の分散サイトストレージユニットへファンアウトを再形成する要求を送信し、前記複数の分散サイトストレージユニットへ送信されるデータを、前記複数の分散サイトストレージユニットからの更新されたブロックリストに基づいて決定し、前記複数の分散サイトストレージユニットへ、前記1次サイトストレージユニットのデータを複製するのに十分なデータをコピーするように適合された前記ロジック

20

をさらに備える請求項1に記載のストレージユニット。

## 【請求項 3】

ストレージユニット(202)であって、

少なくとも部分的に論理ユニット(lun)として構成された前記ストレージ(204)と、

30

前記1次サイトストレージユニットとのファンアウト関係が存在する論理ユニット(lun)を対象にした書き込みを検出し、ストレージ論理ユニット(lun)の変更されたブロックを追跡し、前記複数の分散サイトストレージユニットへ非同期通信パケットを送信するように適合された前記ロジックと、

1つ又は複数の書き込みによって変更されたブロック番号と、ブロックの内容の一意の識別子を示すシーケンス番号とを含むデータパケットを収集し、前記複数の分散サイトストレージユニットへ非同期通信によって前記データパケットを送信するように適合された前記ロジックと、

複数のデータパケットをパケットグループに結合し、前記複数の分散サイトストレージユニットへ非同期通信によって前記データパケットを送信するように適合された前記ロジックと、

40

前記複数の分散サイトストレージユニットから複数のデータパケット及びシーケンス番号を受信し、前記複数の分散サイトストレージユニット間でデータの内容の相違を求めるように適合された前記ロジックと

をさらに備える請求項1に記載のストレージユニット。

## 【請求項 4】

ストレージユニット(202)であって、

データの再形成を開始する信号を受信し、前の複製書き込みストリームから保留中のバッファリングされた書き込みの処理を終了し、1次サイトとして新たに動作するストレージユニットへ、前記複製によって更新されたブロックのリストを含む、更新されたブロッ

50

クリストを送信するように適合された前記ロジック  
をさらに備える請求項 1 に記載のストレージユニット。

【請求項 5】

ストレージユニット ( 2 0 2 ) であって、  
フェイルオーバーするコマンドを受信し、2 次サイトストレージユニットとしてのオペレーションから新たな 1 次サイトストレージユニットに再構成し、前記複数の分散サイトストレージユニットの残りの 2 次サイトストレージユニットに、ファンアウトが再形成されていることを通知する信号を送信するように適合された前記ロジックと、  
新たな 1 次サイトストレージユニット用に動作可能な前記ロジックであって、  
前記複数の分散サイトストレージユニットの残りの 2 次サイトストレージユニットが  
前記新たな情報源よりも新しい状態を有するかどうかを判断し、  
前記より新しい状態を有する前記 2 次サイトストレージユニットへ、前記新たな 1 次  
サイトストレージユニットにおいて新しくないデータの要求を送信し、  
前記要求されたデータで前記新たな 1 次サイトストレージユニットを更新する、  
ように適合された  
前記ロジックと  
をさらに備える請求項 1 に記載のストレージユニット。

【請求項 6】

ストレージシステム ( 1 0 0 ) であって、  
1 : n ファンアウト構成で配置された複数のストレージアレイ ( 1 0 2 ) と、  
前記複数のストレージアレイ ( 1 0 2 ) において実行可能なロジック ( 1 0 4 ) であって、情報源ストレージアレイ ( 1 0 2 S ) から n 個の宛先ストレージアレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) へのデータ複製中のデータの変更を追跡し、追跡されたデータ変更情報を前記 n 個の宛先ストレージアレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) 間で相互に共有し、1 : n - 1 ファンアウト構成に再形成することによってフェイルオーバー状態に対応するように適合されたロジックであり、前記再形成は、前記 n 個の宛先ストレージアレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) からの前記相互に共有された、追跡されたデータ変更情報に従って指示されるロジックと  
を備えるストレージシステム。

【請求項 7】

ストレージシステム ( 1 0 0 ) であって、  
前記障害を受けた情報源ストレージ ( 1 0 2 S ) を除外するように前記複数のストレージアレイ ( 1 0 2 ) を構成することにより前記フェイルオーバー状態に対応し、前記フェイルオーバーとほぼ同時に行われる割り当てにおいて、前記 n 個の宛先ストレージアレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) の 1 つを新たな情報源ストレージアレイ ( 1 0 2 S ) として動作するように割り当てるように適合された前記ロジック ( 1 0 4 ) と、  
残りの n - 1 個の宛先ストレージアレイのデータを前記新たな情報源ストレージアレイに準拠して再形成することにより、前記フェイルオーバー状態にさらに対応するように適合された前記ロジック ( 1 0 4 ) と  
をさらに備える請求項 6 に記載のストレージシステム。

【請求項 8】

ストレージシステム ( 1 0 0 ) であって、  
個々の宛先ストレージアレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) において実行可能なロジック ( 1 0 4 ) であって、前記情報源ストレージアレイ ( 1 0 2 S ) とのファンアウト関係が存在する論理ユニット ( l u n ) を対象にした書き込みを検出するように適合され、他の前記複数の宛先ストレージアレイの複数へ非同期通信パケットを送信するように適合された、ロジック ( 1 0 4 )  
をさらに備える請求項 6 に記載のストレージシステム。

【請求項 9】

ストレージシステム ( 1 0 0 ) であって、

個々の宛先ストレージレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) において実行可能なロジック ( 1 0 4 ) であって、

前記情報源ストレージレイ ( 1 0 2 S ) とのファンアウト関係が存在する論理ユニット ( l u n ) を対象にした書き込みを検出し、

1 つ又は複数の書き込みによって変更されたブロック番号と、ブロックの内容の一意の識別子を示すシーケンス番号とを含むデータパケットを収集し、

他の前記複数の宛先ストレージレイの複数への非同期通信によって前記データパケットを送信する、

ように適合されたロジック ( 1 0 4 ) と、

個々の宛先ストレージレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) において実行可能なロジック ( 1 0 4 ) であって、複数のデータパケットをパケットグループに結合し、他の前記複数の宛先ストレージレイの複数への非同期通信によって前記パケットグループを送信するようにさらに適合されたロジック ( 1 0 4 ) と、

個々の宛先ストレージレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) において実行可能なロジック ( 1 0 4 ) であって、他の前記複数の宛先ストレージレイの複数から複数のデータパケット及びシーケンス番号を受信し、前記複数の宛先ストレージレイ間でデータの内容の相違を求めるように適合されたロジック ( 1 0 4 ) と

をさらに備える請求項 6 に記載のストレージシステム。

#### 【請求項 1 0】

ストレージシステム ( 1 0 0 ) であって、

個々の宛先ストレージレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) において実行可能なロジック ( 1 0 4 ) であって、

フェイルオーバーするコマンドを受信し、

新たな情報源ストレージレイとして再構成し、

前記複数の宛先ストレージレイの残りの宛先ストレージレイに、ファンアウトが再形成されていることを通知する信号を送信する

ようにさらに適合されたロジック ( 1 0 4 ) と、

個々の宛先ストレージレイ ( 1 0 2 D 1、1 0 2 D 2、1 0 2 D 3 ) において実行可能なロジック ( 1 0 4 ) であって、

ファンアウトが再形成されていることを通知する前記信号を受信し、

前の複製書き込みストリームから保留中のバッファリングされた書き込みの処理を終了し、

前記新たな情報源ストレージレイへ、前記宛先の更新されたブロックリストであって、前記複製によって更新されたブロックのリストを含む、更新されたブロックリスト、を送信する

ように適合されたロジック ( 1 0 4 ) と、

前記新たな情報源ストレージレイにおいて実行可能なロジック ( 1 0 4 ) であって、

前記複数の宛先ストレージレイから受信された、更新されたブロックリストの相違を求め、

前記複数の宛先ストレージレイを同期させるのに十分なデータを前記複数の宛先ストレージレイにコピーする

ように適合された、ロジック ( 1 0 4 ) と、

前記新たな情報源ストレージレイにおいて実行可能なロジック ( 1 0 4 ) であって、

前記複数の宛先ストレージレイの或る宛先ストレージレイが、前記新たな情報源ストレージレイよりも新しい状態を有するかどうかを判断し、

最新の状態を有する宛先ストレージレイへ、前記宛先ストレージレイに存在して前記新たな情報源ストレージレイに存在しないデータの要求を送信し、

前記要求に回答して受信されたデータで、前記新たな情報源ストレージレイを更新する

ように適合されたロジック ( 1 0 4 ) と

10

20

30

40

50

をさらに備える請求項 6 に記載のストレージシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ追跡を備えるストレージ複製システムに関する。

【背景技術】

【0002】

データの複数のコピーを保持することは、データが利用可能でない場合、損傷した場合、又は喪失した場合のデータ処理オペレーションのセキュリティ機能の一部である。データ処理システムの組織ユーザ (institutional user) は、一般に、多数の非常に重要な情報を保持し、災害又は惨事に起因して利用不能になることからデータを保護するのに多くの時間及び費用を費やす。冗長なデータコピーを保持するための技法の1つの分類は、ミラーリングと呼ばれる。このミラーリングでは、データ処理システムのユーザは、現場の貴重な情報のコピーを、ローカル又はリモートに位置するリムーバブルストレージ媒体又は2次ミラストレージサイトに保持する。離れた場所ではあるが、たとえば約200キロメートルまでの大都市圏内のリモートミラーリングは、火災、停電、又は盗難を含むローカルな災害からの保護に役立つ。数百キロメートルの地理的距離を超えるリモートミラーリングは、地震、竜巻、ハリケーン、洪水等の惨事からの保護に有益である。多くのデータ処理システムは、複数の地理的距離に位置する複数レベルの冗長性を使用してデータを保護する。

10

20

【0003】

複数サイトのデータ複製及びミラーリングの技術の一態様は、サイトの1つにおける障害及び災害の状態に対応することである。通常、残存しているさまざまなサイト又はノードを同期状態にするために、或るデータの再正規化又は再調整 (reconciliation) が必要とされる場合があり、これは、通常、残存しているノードで再正規化される論理ユニット (lun、複数可) の全面的なコピーを伴うプロセスである。コピーによって、企業クラスの高可用性で且つ災害耐性のあるアプリケーションには許容できない性能及び可用性の劣化が生じる。

【発明の開示】

【課題を解決するための手段】

30

【0004】

ファンアウト関係を再形成するための技法の一実施の形態によれば、ストレージ複製方法は、複数の宛先の中の情報源からデータを複製すること、及び、それら宛先のデータの変更を追跡することを含む。変更の識別情報は、複数の宛先アレイ間で相互に通信される。情報源のフェイルオーバーイベントでは、選択された宛先が新たな情報源として設定され、残りの宛先における複製されたデータを、当該新たな情報源と同期した状態に再形成する。

【0005】

オペレーションの構造及び方法の双方に関するこの発明の実施の形態は、以下の説明及び添付図面を参照することによって最もよく理解することができる。

40

【発明を実施するための最良の形態】

【0006】

ほとんどの障害のシナリオにおいてすべての再正規化のオーバーヘッドをほぼ回避する1:nファンアウトの複数サイトの複製を再構築するためのストレージシステム、ストレージユニット、及び関連した操作技法を説明する。

【0007】

本明細書に開示する技法及び構造を使用してファンアウトを再形成することにより、サイト間トラフィック、再同期時間、及びホストアプリケーションへの性能の影響を削減することができるか、又は、最小にすることができる。この技法及び構造によって、さらに、情報源論理ユニット (lun; logical unit) が障害イベント後に少なくとも1つの対

50

応する、同期されたコピーにアクセスできない期間であるタイムウィンドウを削減することができるか、又は、最小にすることができる。

【0008】

図1A、図1B、及び図1Cを参照して、概略ブロック図が、ストレージシステム100の一実施の形態を示している。このストレージシステム100は、1:nのファンアウト構成、例示では1:3のファンアウト構成に配列された複数のストレージレイ102S、102D1、102D2、及び102D3を備える。図1Aは、1つの情報元レイ102S及び3つの宛先レイ102D1、102D2、及び102D3を有する1:3の論理ユニット(1un)ファンアウトの例を示している。情報元レイは、ハブと呼ばれることがある。図1Bは、情報元レイ102Sの障害時の1:31unファンアウトの例を示している。図1Cは、情報元の障害後であって、且つ、ファンアウトを1:2構成に修復した後のストレージシステム100の構成を示している。ロジック104は、複数のストレージレイ102S、102D1、102D2、及び102D3にわたって分散され、これらのストレージレイにおいて実行可能である。いくつかの構成では、このロジックは、ストレージレイの外部のホスト、コンピュータ、コントローラ、ストレージ管理デバイス等にまで及ぶことがある。ロジック104は、情報元ストレージレイ102Sからn個の宛先ストレージレイ102D1、102D2、及び102D3へのデータの複製中におけるデータの変更を追跡するように適合されている。ロジック104は、追跡されたデータ変更情報をn個の宛先ストレージレイ間で経路106を介して相互に共有する。ここで、nは任意の適切な数である。ロジック104は、1:n-1のファンアウト構成に再形成することによってファイルオーバ状態に対応する。再形成は、相互に共有された、n個の宛先ストレージレイ102D1、102D2、及び102D3からの追跡されたデータ変更情報に基づいて指示される。

10

20

【0009】

簡単なリモート複製デプロイメント(deployment)は、1:1構成とも呼ばれる2つのサイトとすることができる。この1:1構成では、1つの論理ユニット(1un)に対する入出力オペレーションが、通常は別々の地理的位置の宛先レイにおける宛先1unに対してリアルタイム(real time)で複製される。イベント、たとえば、天候、地震、停電、破壊状況等の災害状態が1次サイトに影響を与える場合、アプリケーションは、2次サイトにオペレーションを移動させることによって、そのイベントの直前の状態に回復することができる。1:1のリモート複製配置の制限は、サイトのイベントに続いて、損傷を受けたサイトが回復されるまで、複製されたデータの単一のコピーしか残っていないということである。回復時間は、かなりかかる場合があり、銀行業、仲買業、株式取引、軍隊、健康管理等の産業及び機関における、災害耐性があり、且つ、高可用性のアプリケーションの要求に対して許容できない単一障害点のリスクを表す。多くの、災害耐性があり、且つ、高可用性のユーザは、1つのサイトがイベントによって除かれる場合、2つのアクティブなサイトとなる3サイト複製の仕様を課す。

30

【0010】

論理ユニット(1un)のファンアウトは、情報源1unのミラーコピーを2つ又は3つ以上の宛先レイに同時に複製することを伴うレイベースのリモートアプリケーション技術である。その情報源への新たな書き込みは、非同期又は同期のいずれかの方法で、複数の宛先へ複製される。同期複製では、情報源1unに対する書き込みオペレーションは、その書き込みが情報源1un及び宛先1unの双方にコミットされた時に完了したものとして開始ホストに通知される。非同期複製では、情報源1unに対する書き込みオペレーションは、その書き込みが情報源1unへはコミットされたが、宛先1unへはコミットされていない時に完了したものとして開始ホストに通知される。書き込みは、情報源1unに対する書き込みとは独立した動作で後に宛先1unに適用される。非同期複製によって、情報源1unに対する書き込みオペレーションに書き込み待ち時間の遅延を被らないので、地理的に分散したリモート複製の性能を最高レベルにすることが可能になる。同期複製は、距離によって性能は低下するが、宛先1unが、常に、バイト単位(byte-w

40

50

ise) で正確な情報源  $lun$  の複製又はほぼ正確な複製であることを保証する。

【0011】

複数のサイトのリモート複製は、単一の  $lun$  ファンアウトである簡単な 1 : 2 ファンアウト技術を使用して実施することもできるし、例示のストレージシステム 100 の 1 : 3 ファンアウトのように実施することもできる。

【0012】

$lun$  ファンアウト比が高くなるほど、冗長性が増加し、したがって信頼性が増加する。また、 $lun$  ファンアウトによって、ユーザに対するデータのアクセス性も可能になる。たとえば、ブロードキャスト分散データの分散モデルは、 $n$  が 2 以上、又は、それよりもはるかに大きい 1 :  $n$  ファンアウトを備えることができる。特定の例では、ライブストリーミングメディアの供給は、クライアントアプリケーションの近くにあるサーバに適用することができ、それによって、大きなネットワークのオーバーヘッドが取り除かれる。

10

【0013】

ストレージシステム 100 は、特定の位置、たとえばロンドン等の地理的位置の情報源ストレージアレイ 102 S 又はハブで視覚化 (visualize) することができる。ハブ 102 S から外部には、通信リンク 108 が伸びている。この通信リンク 108 は、ハブ 102 S をリモートストレージアレイ 102 D 1、102 D 2、及び 102 D 3 に接続する。ハブ 102 S は、情報源論理ユニット ( $lun$ ) 110 S を含むアレイとすることができる。リモートストレージアレイ 102 D 1、102 D 2、及び 102 D 3 は、リモート  $lun$  110 D 1、110 D 2、及び 110 D 3 を含む。データは、通信リンク 108 上を同期又は非同期のいずれかで流れる。一般的な場合、ストレージアレイは地理的に分散している。単なる例として、第 1 の宛先アレイ 102 D 1 及び第 1 の宛先  $lun$  110 D 1 はニューヨークに配置することができ、第 2 の宛先アレイ 102 D 2 及び第 2 の宛先  $lun$  110 D 2 は東京に配置することができ、第 3 の宛先アレイ 102 D 3 及び第 3 の宛先  $lun$  110 D 3 は香港に配置することができる。幅広い分散によって、限られた地理的領域で発生し得る障害の回避が容易になる。通常の構成では、1 つ又は複数のリンクは非常に遠隔で、且つ、非同期であり、1 つのリンクは大都市圏内にあるか又は大都市圏全体にわたるものであり、且つ、同期したものであり、情報源  $lun$  は、応答性をかなり良くすることができると同時に、適切な災害耐性を維持することができる。他の構成も可能である。

20

30

【0014】

図 1 B は、1 :  $n$  ファンアウトオペレーションに関する問題を示し、情報源 102 S の喪失に起因するファンアウト関係の破壊、ファンアウト関係の再形成、残っているストレージアレイ 102 D 1、102 D 2、及び 102 D 3 でファンアウト複製の再確立に対応する (address) 方法を示している。

【0015】

情報源ストレージアレイ 102 S 又はハブが喪失する状態又はハブへの通信が喪失する状態が発生した場合に、ストレージシステム環境 100 が宛先ストレージアレイ 102 D 1、102 D 2、及び 102 D 3 の 1 つにオペレーションをフェイルオーバできるときは、アプリケーションを続行することができる。

40

【0016】

図 8 A に示すように、追跡情報の追跡及び共有を含まないファンアウト配置 800 では、宛先アレイ間に関連は存在しない。各宛先アレイ 802 D は、ハブ 802 S との関係のみを有する。図 8 B に示すように、ハブ 802 S の喪失によって、ファンアウト関係が再形成される時、宛先アレイ 802 D は、どのブロックが、他の宛先アレイ 802 D の対応する  $lun$  810 D に書き込まれたのか、又は、書き込まれていないのかに関する情報を有しない。その結果、新たなハブとして動作するように決定された宛先アレイは、 $lun$  を全面的にコピーしなければならず、これは、かなりの時間、おそらく数日続くおそれのあるオペレーションであり、この全面的なコピーオペレーションによって被る性能ペナルティは、重大なものとなる可能性がある。ファンアウトの形成が一旦開始されると、顧客

50

は、ファンアウトの開始の程度にかかわらず、データの単一の良好なコピーのみが保護される状況に置かれる。

【0017】

図8A及び図8Bは、以下でより詳細に説明する。

【0018】

図1Aに示すストレージシステム100を再び参照して、データ保護及び効率性は、宛先ストレージアレイ102D1、102D2、及び102D3の情報間の現行の関連を維持することによって高められる。この関連は、各宛先アレイが、他の宛先ストレージアレイ102D1、102D2、及び102D3のいずれかにおけるパートナーとなるlun110D1、110D2、及び110D3に関して異なる、ローカルlun110Sのブロックの識別情報及び追跡を保持することを可能にする技法のオペレーションを通じて維持される。この技法はアカウントング技法 (accounting technique) と呼ばれることがある。

10

【0019】

通常、個々の宛先ストレージアレイ102D1、102D2、及び102D3は、各宛先論理ユニット (lun) 110D1、110D2、及び110D3のデータブロックの変更を追跡するように構成されたロジックを含む。追跡オペレーションにおいて、このロジックは、情報源ストレージアレイ102とのファンアウト関係が存在する論理ユニット (lun) を対象にした書き込みを検出することができ、他の宛先ストレージアレイ102D1、102D2、及び102D3のそれぞれへ通信パケットを送信することによってその書き込みオペレーションに回答することができる。例示の実施の形態では、この通信パケットは、相互リモート通信リンク106において宛先ストレージアレイ102D1、102D2、及び102D3間で交換される。いくつかの実施の形態では、宛先ストレージアレイ102D1、102D2、及び102D3は、非同期通信によって通信し、それによって、通信の応答が進行することを待つことなくネットワークにおいて要求が行われる。この応答は、後に来ることがある。

20

【0020】

例示の実施の形態では、ロジックは、1つ又は複数の書き込みによって変更されたブロック番号と、情報源アレイ102Sによって供給されて情報源102Sに対する書き込みを示すシーケンス番号とを含むデータパケットを収集する。

30

【0021】

情報源アレイ102Sでは、110S等のファンアウトを受ける論理ユニット (lun) について、ブロック書き込みは、シーケンス番号を割り当てることによってハンドリングされる。各書き込みは、通常、1つのシーケンス番号によって特定される。情報源アレイ102Sは、ブロック識別子 (ID) と、そのブロックのデータと、シーケンス番号とを含む情報を通信リンク108ですべての宛先ストレージアレイ102D1、102D2、及び102D3へ送信する。

【0022】

ブロックは、物理エンティティとすることもできるし、論理エンティティとすることもできる。たとえば、ブロックはトラック/セクタとすることができる。このトラック/セクタは、通例、物理的な記憶素子又はディスクに関係する。一般的な論理エンティティは、lunにおける連続したブロック番号である。一般に、ブロックは、記憶素子の定義されたオフセットにおける既知の一定サイズのデータチャンクである。異なるタイプのストレージアレイは、異なるタイプのブロックを使用してデータを複製することができる。たとえば、いくつかのアレイはトラック/セクタアイテムとして指定された物理ブロックを使用し、他のアレイは論理ブロックを使用する。

40

【0023】

例示の目的で、且つ、論理ブロックの説明例として、lunは、1000個のブロックを含むことができ、各ブロックは、1メガバイト (MB) のサイズを有する。この状況におけるブロック562の意味は、lunにおける562番目の1MBのブロックというこ

50



とである。

【0024】

シーケンス番号は、情報源アレイによって割り当てられる。シーケンス番号は、特定の論理ユニット ( lun ) について連続的に順次インクリメントする整数である。情報源アレイにおける第1の論理ユニット、たとえば lun A に対する書き込みは、同じ情報源アレイにおける lun B 等の第2の論理ユニットのシーケンス番号には影響を与えない。したがって、シーケンス番号は、その情報源 lun の各書き込みオペレーションの到着ごとに1つずつインクリメントする。読み出しオペレーションでは、シーケンス番号は変更されずにそのままにされる。

【0025】

パケットは、効率性を促進するために、パケットグループに束ねて、宛先ストレージアレイ 102D1、102D2、及び102D3間で通信することができる。宛先ストレージアレイ 102D1、102D2、及び102D3は、さらに、他の宛先ストレージアレイ 102D1、102D2、及び102D3からのデータパケット及び/又はデータパケットグループを相互に受信するように適合され、且つ、そのパケットの情報に基づいてデータの内容の相違を判断するように適合されたロジックも含む。宛先ストレージアレイ 102D1、102D2、及び102D3間のリモート通信リンク106によって、各宛先ストレージアレイは、アレイ間の伝送遅延の影響を受けることを条件として、ほぼすべての時刻においてすべての宛先アレイ間の lun の内容の相違に関する情報を有することが可能になる。リモート通信リンク106は、宛先ストレージアレイ 102D1、102D2、及び102D3間の直接接続とすることができる。いくつかのインプリメンテーションでは、リモートリンク106は、情報源ストレージアレイ 102S への相互接続経路とは独立したものとすることができる。場合によっては、通信リンク106は、同じネットワーク内とすることができる、したがって独立していないものとすることができる。ただし、情報源の近くのリンク106の部分が障害になった場合、たとえば公共のインターネットが動作するように、宛先ストレージアレイ間のリンク106のサブセクション ( subsection ) が使用可能な状態で残っている限り、オペレーションは続行する。リモート通信リンク106は、適切な相互接続技術によって形成することができる。一例は、インターネットプロトコル ( IP ) 通信である。

【0026】

たとえば、図1Bに示すように、ハブアレイが喪失し、図1Cに示すように、ファンアウトが修復されると、宛先ストレージアレイ 102D1、102D2、及び102D3のいずれかを、新たな中央ハブ 102S' として動作するように選択することができる。新たに指定された情報源アレイ又はハブ 102S' は、制御エンティティ、たとえばシステム管理エンティティから、オペレーションをフェイルオーバーするコマンドを受信する。ストレージシステム100のロジックは、フェイルオーバー状態に対応するように適合することができる。この適合は、障害を受けた情報源ストレージアレイ 102S を除外するように複数のストレージアレイを構成し、且つ、フェイルオーバーとほぼ同時に行われる割り当てにおいて、n個の宛先ストレージアレイ 102D1、102D2、及び102D3の1つを新たな情報源ストレージアレイ 102S' として動作するように割り当てるよう複数のストレージアレイを構成することによって行われる。たとえば、新たな情報源の決定は、誘発イベントのタイプ又は性質、日時、テクニカルサポート ( technical support ) の可用性、さまざまなサイトの技術的特徴、さまざまな商慣習等のさまざまな状態又は状況に基づいて行うことができる。たとえば、情報源は、イベントが或る位置では勤務時間中に発生し、別の位置では夜中に発生することに基づいて選択することができる。

【0027】

新たな情報源ストレージアレイ 102S' が一旦選択されると、ストレージシステム100は、残りの n - 1 個の宛先ストレージアレイのデータを、新たな情報源ストレージアレイ 102S' と同期した状態又は準拠した状態に再形成することによってフェイルオーバー状態にさらに対応することができる。ストレージシステム100の残りのものを再形成

10

20

30

40

50

するために、ファンアウトが再形成していることを残りの宛先ストレージアレイに通知するコマンド又は信号を新たな情報源ストレージアレイ102S'から送信することができる。新たな構成における宛先ストレージアレイ102D1'及び102D2'は、この再形成コマンドを受信すると、もはや、前の情報源ストレージアレイ102Sからの新たな要求を受け取ることにはない。肯定応答を元の情報源102Sへ返信できるかどうかにかかわらず、再形成コマンド前に到着したあらゆる書き込みは完了される。また、宛先ストレージアレイ102D1'及び102D2'は、更新を受信した宛先1un110D1'又は110D2'におけるブロックを指定した最終リストを送信することによって、新たな情報源ストレージアレイ102S'に回答する。

#### 【0028】

新たな情報源ストレージアレイ102S'は、残りの宛先ストレージアレイ102D1'及び102D2'から受信した、更新されたブロックリストの相違を判断することによってシステムデータの状態を解明し、ストレージアレイを同期させるのに十分なデータを宛先ストレージアレイ102D1'及び102D2'にコピーバックする。新たな情報源アレイ102S'は、相違するブロックのみを、通信リンク108を介して、再形成している宛先ストレージアレイ102D1'及び102D2'にコピーする。解明段階の期間中、新たな情報源ストレージアレイ102S'は、情報源と比較して宛先アレイで不足しているデータブロックのみを送信し、それら宛先を新たな情報源と同期した状態にする。

#### 【0029】

再形成コマンドが成功する一般的な場合に、新たな情報源ストレージアレイ102S'は、宛先1un110D1'及び110D2'を新たな情報源ストレージアレイ102S'と同期した状態にすることを可能にするために、宛先アレイのそれぞれにどのブロックをコピーするかに関する情報を要求して受信する。新たな情報源ストレージアレイ102S'は、相違するブロックのみを各宛先1un110D1'又は110D2'にコピーし、ファンアウトが再形成される。

#### 【0030】

再形成コマンドを実行できず、且つ、宛先アレイにアクセスできない異常な場合には、アクセスできない宛先アレイはファンアウトの再形成に参加しない。ファンアウトが修復された後、ホストアプリケーションからの任意の書き込みが新たな情報源1un110S'へ送信された場合、説明した技法の一実施の形態を使用すると、アクセスできない宛先アレイがオンラインに戻った時にその宛先アレイを修復してファンアウトに戻すのに、ブロックコピーが使用される。例示の技法を使用すると、当初ファンアウト関係で同期されていた1unが、アクセスできない期間の後にファンに再び参加する状態を得るために、相違するブロックのみをコピーして、最大効率でファンアウトを再同期させることができる。すべての1unブロックの全面的なコピーは、完全に新しい1unがファンに参加する場合にのみ是認される。また、この技法は、元の情報源が1:n-1ファンに再び参加して1:nファンアウトを再形成する場合も包含する。さらに、この技法は、ファンの1つ又は複数の宛先1unがアクセスできない間、書き込みオペレーションを調べる新たな情報源1unの場合も包含する。すべての場合に、この技法は、相違するブロックのみをコピーする動作を含む。

#### 【0031】

また、元の情報源102Sは、1unに適用された書き込みのブロック/シーケンステーブルを保持することもできる。説明した基本的なブロックの相違のアカウントリング及び更新は、書き込みオペレーションに適している。再形成されたリンクが使用可能になるか、又は、元の情報源102Sが宛先として再形成された後にファンに再び参加すると、再形成コマンドに対する応答は、一旦受信されると、通信喪失の継続時間に応じて、比較的大きくなる場合があるブロック/シーケンス番号リストを含むことができる。

#### 【0032】

一般的な実施の形態では、新たな情報源は、1unの複製の完全性にこだわらずに選択される。選択された新たな情報源は、宛先アレイの1つ又は複数と同等の最新の複製を有

10

20

30

40

50

しない場合がある。一般的な実施の形態では、より完全な新たな情報源を得る試みは行われない。しかしながら、いくつかの実施の形態では、最新の情報が求められる場合がある。このような実施の形態は、新たな情報源ストレージレイ 102S' で実行可能なロジックを含むことができる。このロジックは、宛先ストレージレイ 102D1' 又は 102D2' が新たな情報源ストレージレイ 102S' よりも新しい状態を有するかどうかを判断するように適合されたものである。より新しい状態を有する宛先ストレージレイ 102D1' 又は 102D2' は、再形成コマンドの発行及び応答の収集の後に判断される。lunのすべてのブロックにわたって最も高いブロックシーケンス番号を有する宛先ストレージレイ 102D1' 又は 102D2' が最新のものである。いずれかの宛先への通信が遮断された場合、その宛先は交渉に参加することができない。新たな情報源ストレージレイ 102S' は、最新の状態を有する宛先ストレージレイへ、その宛先ストレージレイに含まれるが、新たな情報源ストレージレイ 102S' には存在しないデータの要求を送信する。その宛先から要求されたブロックは、より高いシーケンス番号を有するあらゆるものである。新たな情報源ストレージレイ 102S' は、選択された宛先レイから lunのより新しいブロックを収集し、要求に応答して受信されたデータで、新たな情報源ストレージレイ 102S' を更新する。

10

**【0033】**

新たな情報源レイは、選択された宛先レイから、(lunの)より新しいブロックを収集する。

**【0034】**

図2を参照して、概略ブロック図が、冗長データストレージシステム200での使用に適合されたストレージユニット202の一実施の形態を示している。データストレージシステム200は、少数又は多数のストレージユニット202を有することができる。これらのストレージユニット202は、ストレージ(揮発性メモリ、不揮発性メモリ)204、複数の分散サイトストレージユニットと通信するように適合された入出力(I/O)インターフェース206、及びコントローラ(CPU)208を備える。ファンアウト配置において2次サイトとしてストレージユニット202を動作させ、且つ、1次サイトストレージユニットからストレージ204へデータを複製するように適合されたロジックがコントローラ208において実行可能である。このロジックは、さらに、ストレージ204に書き込まれたデータの変更を追跡し、追跡された変更を複数の分散サイトストレージユニット間で通信する。また、このロジックは、複数の分散サイトストレージユニットから受信された、追跡された変更を収集する。

20

30

**【0035】**

ストレージ204は、ディスクレイ、光ディスクストレージ装置、階層化メモリ、及び/又は、分散しているが密接に制御される、記憶機能(capability)を備えるネットワーク等の任意の適切なストレージ媒体デバイスとすることができる。ストレージ204は、少なくとも部分的に論理ユニット(lun)210として構成される。

**【0036】**

ストレージユニット202が2次サイトストレージユニットとして動作している間、コントローラ208において実行可能なロジックは、1次サイトストレージユニットとのファンアウト関係が存在する論理ユニット(lun)を対象にした書き込みを検出する。このロジックは、書き込みオペレーションによって書き込まれた、論理ユニット(lun)のブロックを追跡する。特定の実施の形態では、この追跡動作は、ストレージユニット202への書き込みによって変更されたブロック番号、及び、ホスト又は情報源から送信された、ブロックの内容の一意の識別子を示すシーケンス番号を収集することを含むことができる。収集されたブロック番号及びシーケンス番号は、データパケットに記憶することもできるし、選択された時間にわたって累積してパケットグループを形成することもできる。パケットグループは、「チャンク」と呼ばれることがあり、たとえば非同期通信により他の分散2次サイトストレージユニット間で直接通信され、追跡された情報が共有される。このロジックは、他の2次ストレージユニットからのパケット及び/又はパケットグ

40

50

ループのブロック番号データ及びシーケンス番号データを受信し、ストレージユニット 202 にローカルな情報に関する情報を解析して、通常は分散サイトにおいて、複数の 2 次ストレージユニット間のデータの内容の相違を求める。

#### 【0037】

ストレージユニット 202 は、フェイルオーバーするコマンドを受信することができる。このコマンドに回答して、コントローラ 208 において実行可能なロジックは、ストレージユニット 202 を 1 次サイトとして動作させ、ストレージユニット 202 にネットワーク接続された分散サイトストレージユニットへファンアウトを再形成する要求を送信する。分散サイトストレージユニットは、リモートユニットにおいて複製された書き込みを示す、更新されたブロックリストを送信することによって、この再形成する要求に回答する。ストレージユニット 202 は、更新されたブロックリストに基づいて、どのデータを分散サイトストレージユニットへ送信するかを判断し、そのデータを分散サイトユニットにコピーする。コピーされたデータは、1 次サイトからのデータ ( l u n ) の正確なバイト単位の複製を作成するのに十分なものである。

#### 【0038】

フェイルオーバー後、ストレージユニット 202 が、1 次として動作するように選択されない場合、コントローラ 208 において動作しているロジックは、通常はフェイルオーバーの最初の表示として、新たな 1 次からの、データの再形成を開始する信号又はコマンドを受信する。この再形成する信号に回答して、ロジックは、前の 1 次からの新たな要求の受け取りを直ちに停止し、新たな 1 次へ、元の複製ストリームによって更新された最後のブロックのリストを含む更新されたブロックリストを送信する。

#### 【0039】

図 3 を参照して、概略フローチャートは、危険性なくデータを維持しつつそのデータを複数の地理的位置に複製するように、ファンアウト関係を高速に再形成するよう適合された技法の一実施の形態を示している。ストレージ複製方法 300 は、情報源からのデータを複数の宛先間で複製すること ( 302 ) と、複数の宛先におけるデータ変更を追跡すること ( 304 ) とを含む。変更の識別情報は、複数の宛先アレイ間で相互に通信される ( 306 )。情報源のフェイルオーバーイベント 308 において、選択された宛先が新たな情報源として設定され ( 310 )、残りの宛先の複製データは、新たな情報源と同期した状態に再形成される ( 312 )。通常、選択された宛先は、ユーザがグラフィカルユーザインターフェース ( G U I ) のボタンを押したり、コマンドラインインターフェース ( C L I ) にコマンドをタイプ入力したりしてフェイルオーバーを起動すること等の動作により、新たな情報源として設定することができる ( 310 )。

#### 【0040】

フェイルオーバー 308 の前における宛先間の変更データの相互通信 306 によって、ハブアレイを伴う障害後に、データの現在のコピーが 1 つしか存在しない状態にユーザアプリケーションが置かれる時間を大幅に減少させることが可能である。また、変更データの通信 306 によって、全面的なデータコピーが回避されるので、フェイルオーバー後の情報源 l u n に対するスループット性能が改善される。

#### 【0041】

図 4 を参照して、概略図が、ストレージシステムにおけるデータ追跡の一例を示すブロックマップシーケンス 400 を示している。ファンアウトオペレーションの開始時において、1 つ又は複数の宛先がファンアウトに追加されると、データの全面的なコピーが通信リンクによって宛先に送信され、情報源及び宛先における対応する論理ユニット ( l u n ) のデータが同期される。ブロック識別子 ( I D ) 404 及びシーケンス番号 406 の完備したブロックマップ 402 が、開始時に情報源アレイから宛先アレイへ送信される。全面的なコピーが、情報源から宛先への通信経路で完了され、宛先は、情報源 l u n に対応する l u n のブロック番号及びシーケンス番号の完備したマップを有するようになる。l u n のすべてのブロックは、2 列の値の配列 402 で表される。第 1 列 404 はブロック番号である。第 2 列 406 は、ブロック列 404 の各ブロックに関連したシーケンス番号

10

20

30

40

50

を含む。どの行も、異なる一意のシーケンス番号を有する。情報源アレイにおいて実行される例示のデータ追跡技法によれば、重複したシーケンス番号は許されず、生じる (occur) ことはできない。その後、情報源アレイが書き込みを受信すると、ブロックID、データ、及びシーケンス番号を含む情報3つ (triplet) が各宛先へ通信される。各宛先におけるこのブロックマップのエントリは、書き込みがコミットされた時に、新たなシーケンス番号で重ね書きされる。

#### 【0042】

例示のために、ブロックマップ402は、情報源アレイで形成され、1つ又は複数の新たな宛先アレイへ通信される5ブロック lun の非常に単純化した例を示している。各宛先は、情報源アレイに記憶されたブロックマップテーブル402に関連したテーブルを保持する。

10

#### 【0043】

いくつかの実施の形態では、情報源が、その後、宛先としてファンアウトに再び参加する場合の使用に備えて、情報源が各宛先の列も保持するように、情報源の追跡テーブルを拡張することができる。情報源は、いかなる場合にも、宛先への書き込みに関する肯定応答を受信するので、追加オーバーヘッドがほとんどないか又は全くなく、それらの列は保持される。情報源におけるすべての宛先ブロックの追跡によって、情報源は、情報源 lun に影響を与えない障害イベントに続いて全面的なコピーを行うことなく、ファンアウトに再び参加することが可能になる。したがって、例示の技法により、1:n-1ファンアウトから1:nファンアウトに戻って再形成することが可能である。

20

#### 【0044】

通常、情報源アレイは、書き込みシーケンスの個々の書き込みとして宛先アレイへ書き込みを送信することができる。いくつかのインプリメンテーションにおいて、又は、いくつかの条件下において、情報源アレイは、複数の書き込みを累積するか、又は、束ねることができる。束ねられた書き込みの通信では、同じブロックが、その束に2つ以上の書き込みを有する場合、最後のシーケンス番号及び関連したデータビットのみが、そのブロックの宛先 lun へ送信される。したがって、特定のブロックが繰り返し書き込まれる状況では、そのブロックのデータは、チャンクが構築される単位時間あたりに1回しかリンクによって送信されないため、束ねることによって、効率性が改善される。宛先 lun がクラッシュコンシステントな状態 (crash-consistent-state) に常にあることを確保するのに、トランザクショナルセマンティクス (transactional semantics) を使用することができる。クラッシュコンシステントな状態では、lun は、チャンクの適用前の正確なバイト単位の値又はチャンクの適用後の正確なバイト単位の値のいずれかを含む。宛先 lun が、部分的なチャンクしか適用されていない状態に入る場合、書き込みオペレーションは、情報源 lun と同じ順序で宛先 lun に適用されていないので、そのチャンクは、おそらくクラッシュコンシステントではない。チャンクデータの移動及びクラッシュコンシステントは、重要な影響を例示の技法にほとんど与えないか、又は、全く与えないが、トランザクショナルセマンティクスによって、どの宛先がファンアウトの新たなハブとして選ばれるかについての判断を容易にすることができる。したがって、チャンク手法によって、データのいくつかのブロック及び対応するシーケンス番号は宛先アレイに決して送信されないようにすることができ、したがって、このような重ね書きされるシーケンス番号は、どの宛先テーブルにも決して現われないようにすることができる。このような省略されたシーケンス番号は、例示の技法の動作性には重要ではない。

30

40

#### 【0045】

情報源アレイから宛先アレイのすべてへの更新に加えて、宛先アレイは、宛先アレイ間の相互接続を介して更新も受信する。宛先アレイ間の相互接続は、ブロック番号及びシーケンス番号の組み合わせの更新も供給する。

#### 【0046】

ブロックマップ408 T1は、時刻 T1における第1の宛先アレイである宛先1のブロ

50

ックマップを示している。このブロックマップ408T1は、情報源アレイのブロックマップ402の対応する列と同様に、宛先1に独自のブロックID410及びシーケンス番号412を含む。これに加えて、ブロックマップ408T1は、相互接続された他の宛先アレイのシーケンス番号も保持し、ここでは、宛先2及び宛先3のシーケンス番号を、各宛先2のシーケンス番号414及び宛先3のシーケンス番号416に保持する。例示では、宛先2のシーケンス番号は、ブロック3についてのみ宛先1と異なる。宛先3のシーケンス番号は、ブロック2及び3について宛先1と異なる。これらの不一致は、アレイ間のさまざまな通信遅延又は書き込みを束ねるために被るアレイのさまざまな内部遅延に起因する場合があります、宛先のアカウントングビュー（accounting view）を遅らせる場合がある。同期複製の場合には、多くとも2～3の不一致しか予想されない。非同期複製の場合には、不一致の発生率は変化し、場合によっては大きくなる可能性がある。例示の技法は、どの宛先lunが別の宛先lunに先行しているかにかかわらず、且つ、宛先lunのいずれかが互いにどれだけ遅れているか又は先行しているかにかかわらず、フェイルオーバー時の不一致を解消する。また、例示の技法は、データ移動を削減するか、又は、最小にする。

10

**【0047】**

各宛先は、適切なlunについて同様のブロックマップテーブルを保持し、更新する。

**【0048】**

たとえば、フェイルオーバーインシデント（incident）の時刻T2において、少なくとも一時的に情報源サイトを除去するイベントを挙げる。この例では、宛先1が新たな情報源アレイとなるように選ばれる。宛先1は、「再形成」コマンドと、宛先1が適用可能なlunの情報源アレイとして制御を行っていることを示す指令とを宛先2及び3へ送信する。宛先2及び宛先3は共に、元の情報源アレイからの新たな書き込みパケットの受け取りを停止し、宛先がコミットした最終的な1組のブロック番号及びシーケンス番号の対で、新たな情報源アレイである前の宛先1に回答する。次に、宛先1は、最終時のマップ408T2として示すブロックマップを更新する。

20

**【0049】**

前の宛先1は、新たな情報源アレイとして、ブロックマップテーブル408T2をスキャンして、一致しない行エントリの検出を可能にする。例示では、宛先2のブロック3並びに宛先3のブロック2及び3は、新たな情報源アレイのエントリと一致しない。したがって、新たな情報源アレイは、すべてのデータビットを含むブロック3の内部コピーを宛先2へ送信し、ブロック2及び3の内部コピーを宛先3へ送信する。新たな情報源から宛先2及び3へのコピーの完了に続いて、宛先2及び3の対応するlunは、前の宛先1と同じ正確なブロックごとの内容を含む。オペレーションは、1：2ファン構成が前の1：3構成に取って代わり、且つ、前の宛先1が新たな情報源アレイとして実行する追跡状態に戻る。データのコピーに続いて、すべてのアレイは同期状態になる。例示では、全面的なデータコピーは、不一致のブロックについてのみ行われ、7つの一致したブロックの全面的コピーは不要とされる。この特定の例では、この技法は、例示のデータ追跡を使用しない技法の30%のコピー負担しか有しない。lunあたり5ブロックよりもはるかに多くのブロックを有する実世界の例では、ほとんどの使用のシナリオについて、この節減は大幅に大きくなり、通常、コピー負担は、すべてのブロックの全面的コピーと比較して10%以下となる。

30

40

**【0050】**

この例は、選択された宛先に対するフェイルオーバーを示している。別の実施の形態のインプリメンテーションでは、フェイルオーバーは、ファンが再形成される時に、選択された宛先がデータの最新のコピーを継承する、一般的に好まれる宛先に対して行うことができる。この技法は、フェイルオーバーが、最新のコピーを含む宛先、ここでは宛先3に対して行われる点を除いて、前に示したのと同じシナリオ及び動作を伴う。ブロックマップテーブル418は、テーブルを完成させるハンドシェイク後の1：2ファンアウト構成の開始状態のステータスを示している。ブロックマップテーブル418は、宛先3の視点からブ

50

ロックマップテーブル408T2を見たものである。新たな情報源アレイとしての宛先3は、テーブルをスキャンし、最も高いシーケンス番号が宛先1に含まれると判断する。したがって、宛先3は、異なるあらゆるブロックを転送するように宛先1に要求する。この例では、宛先3は、ブロック2及び3のデータビットを要求する。転送後、宛先3は最新のデータを有する。宛先3は、例示の方法に従って、宛先3から宛先2にブロック3をコピーすることにより宛先2を宛先3と等しい状態にすることができる。その結果、すべての宛先は、最新のデータを含み、宛先3は、新たな情報源としてのオペレーションを開始する準備ができる。

#### 【0051】

「最新のデータ」の概念は、フェイルオーバー時にアクティブな通信リンクを有する宛先に適用される。実際の最新のデータを含む宛先が、リンク障害のためアクセス不能である場合、新しさでは劣るが、他のどのアクセス可能な宛先よりも新しいデータを有するアクセス可能な宛先が、「最新のデータ」を有するものとみなされる。

#### 【0052】

図5を参照して、概略テーブル図が、複製追跡中にデータを収集するための使用に適したデータ構造体500の別の実施の形態を示している。例示のブロック列は、トラック/セクタデータを指定し、いくつかの実施の形態では、図4に示すブロック識別子を説明するための異なる技法として使用することができる。このデータ構造体は、通常は書き込みオペレーションを介してアクセスされたブロックのテーブルを含むさまざまなファイル、テーブル、サイドファイル等として実施することができる。情報源データ構造体502は、1次ストレージに関連したオブジェクト又はアプリケーションである。たとえば、情報源ハブは、書き込みのテーブルを保持し、テーブルに対する変更を他のアレイ又は宛先へ転送する。データ構造体502Sによって示すように、情報源502は、ホストから書き込みを受信し、それら書き込みを逐次的又はチャンクのいずれかで各宛先に分配する。宛先データ構造体504D1、504D2、及び504D3は、3つの宛先ストレージにそれぞれ関連した、対応するオブジェクト又はアプリケーションである。宛先データ構造体504D1、504D2、及び504D3は、個々の宛先によって受信されたデータを示している。このデータは、まだストレージにコミットされていないものである。宛先は、書き込み、たとえば構造体502Sに示す書き込みのストリームを受信し、それらの書き込みを受信時に順に1unに適用する。チャンク技法が使用され、ブロックの重ね書きが行われた場合、情報源502Sは、単一行と同じブロックに対する上書きを取り込み、書き込みのチャンクは宛先に対する単一ランザクションとして適用される。そうではなく、ストリーミングされるインプリメンテーション又はチャンクでないインプリメンテーションにおいて、書き込みが単一行に対して行われ、ブロックが重ね書きされず、且つ、同じブロックが複数の行に複数回表される可能性がある場合、それらの書き込みは、各宛先において順に適用することができ、ランザクションサイズは単一行となる。各宛先は、ストリーミングされた送信又はチャンク送信のいずれかによって構造体502Sの情報を受信する。各宛先は、情報源からテーブルの変更のリストを受信し、変更されたデータのブロック及びシーケンス番号のリストを他のすべての宛先へ転送する。この変更情報のリストは、リアルタイムで転送することもできるし、蓄積して、選択された蓄積の後に転送することもできる。宛先504D1、504D2、及び504D3が、逐次的な送信又はチャンク送信のいずれかを使用して構造体502Sのデータを受信すると、各ブロック又はチャンクはディスクにコミットされるので、宛先は、1組のブロック識別情報及びシーケンス番号の組み合わせを他の宛先アレイへ送信する。このデータの組み合わせは、コミットされたデータのブロック番号及びシーケンス番号を示す。コミットされたブロックのデータバイトは、或る宛先から別の宛先に送信されず、それによって、宛先間の帯域幅利用が最小限にされて、このアカウントング技法は効率的なものとなる。

#### 【0053】

これらのさまざまなデータ構造体は、各情報源ストレージ及び宛先ストレージのそれぞれにつき、ブロックフィールド506S、506D1、506D2、及び506D3、デ

10

20

30

40

50

ータフィールド508S、508D1、508D2、及び508D3、並びにシーケンスフィールド510S、510D1、510D2、及び510D3を含む。ブロックフィールド506S、506D1、506D2、及び506D3は、書き込みが行われるストレージ媒体の論理ユニット(Lun)やトラック/セクタ情報等の1つ又は複数の位置を指定する。データフィールド508S、508D1、508D2、及び508D3は、対応するブロックフィールドの各トラック/セクタ情報に対応する、バッファリングされたデータを示す。シーケンスフィールド510S、510D1、510D2、及び510D3は、情報源によって定義されたシーケンス番号であって、対応するデータフィールドに列挙された各データ及びブロックフィールドに列挙されたトラック/セクタ情報に関連したシーケンス番号を特定する。

10

**【0054】**

いくつかの実施の形態では、データ構造体は、特定のエントリが他のストレージユニットに関係付けられたことの肯定応答を指定する肯定応答フィールドを含むことができる。たとえば、肯定応答フィールドの論理値1は、特定のシーケンス番号のエントリが他のストレージユニットに複製されたことを示す信号を当該他の2次ストレージユニットから受信したことを示すことができる。論理値0は、特定の2次ストレージユニットに対する複製が行われていないことを示すことができる。

**【0055】**

オペレーションの例では、情報源は、4から9の複製されたシーケンス番号を示している。第1の宛先は、シーケンス番号4から8に対応する書き込みオペレーションを複製する。第2の宛先は、情報源の書き込みのすべてを複製する。第3の宛先は、シーケンス番号4から8に対応する書き込みオペレーションを複製する。異なるストレージユニット間の相違は、リンクの一時的な動作不能に起因することもあるし、同期通信又は非同期通信のいずれかを介して通信できるリンク間のタイミングの相違によることもある。リンク間が非同期通信である場合、多くの書き込みの完了に相違が生じる場合があり、したがって、多くのシーケンス番号に相違が生じる場合がある。リンク間が同期通信である場合、通常、完了時の相違は多くとも1つの書き込みとなる。

20

**【0056】**

フェイルオーバ状態の場合、データは、テーブルのシーケンス番号の識別情報に基づいて、新たな情報源の状態に回復される。例示の技法では、どの資源が特定のブロックに遭遇したかのアカウントिंगを管理するために、データではなくシーケンス番号をストレージユニット間で送信することによって、データトラックは削減される。

30

**【0057】**

図6A及び図6Bを参照して、フローチャートが、情報源の障害の場合に、ファンアウトを削減されたファンアウト比に修復するための技法の実施の形態を示している。ストレージ複製方法600は、情報源から複数の宛先にデータを複製すること(602)、及び、情報源のフェイルオーバ状態を検出すること(604)を含む。フェイルオーバ状態と同時に発生した状態に基づいて、新たな情報源が複数の宛先の中から選択される(606)。この新たな情報源は、複数の宛先でデータの再形成を開始する信号を送信する(608)。

40

**【0058】**

フェイルオーバ時に利用可能な情報及び状態に基づいて取り替える情報源を選択することにより、障害の位置及び原因、レスポンス(response)を実行する資源の可用性、ストレージシステムの各部分の作業負荷(workload)等の因子に基づいて効率的なレスポンスが可能になる。フェイルオーバーイベントが行われるまで、最も適した対応は知られていないので、複数の宛先の中から新たな情報源を同時に選択することによって、柔軟なオペレーションが促進される。

**【0059】**

さらに、この方法は、図6Bに示すように、データ複製中に、個々の宛先の複製ステータス情報を複数の宛先全体にわたって分配すること(610)も含むことができる。この

50



複製ステータス情報を使用すると、複数の宛先において、新たな情報源と同期状態にデータを再形成することができる(612)。新たな情報源で複製ステータス情報が利用可能であることによって、新たな情報源ハブの lun を対象にした入出力オペレーションが、複製の復旧(re-establishment)中に削減されるか、又は、最小にされるので、性能の改善が可能になる。同様に、複製ステータス情報が個々の宛先で利用可能であることによって、宛先 lun を対象にした入出力オペレーションも、複製の復旧中に削減されるか、又は、最小にされるので、性能の改善が可能になる。

#### 【0060】

この改善は、情報源での入出力オペレーションが再形成オペレーションの一部としての大きなボリュームのデータのコピーと競合する必要がなく、宛先での入出力オペレーションも再形成オペレーションの一部としての大きなボリュームのデータのコピーと競合する必要がないことから得られる。同様に、複製の復旧中におけるサイト間リンクの帯域幅に対するこの性能の影響も、削減されるか、又は、最小にされる。この技法によって、情報源が障害になった場合の再形成のための相互通信を制限することが可能になり、帯域幅及び他の資源を大幅に消費する全面的なコピーを回避することが可能になる。すべての場合において、性能の改善は数桁も大きくなる可能性がある。その結果、例示の技法を使用した 1:n ファンアウト技術は、ホスト側のアプリケーションを、性能を劣化させることなく稼働させ続けたい高可用性の災害耐性のあるユーザには非常に魅力的なものとなる可能性がある。

10

#### 【0061】

図7A及び図7Bを参照して、フローチャートが、情報源の障害の発生時にファンアウト構成を再形成するための技法の実施の形態を示している。例示の情報源複製方法700は、情報源から複数の宛先へデータを複製すること(702)、及び、データの再形成を開始する信号を宛先において受信すること(704)を含む。宛先では、前の複製書き込みストリームから保留中の、バッファリングされた書き込みの処理が終了される(706)。宛先は、更新されたブロックリストを新たな情報源へ送信する(708)。更新されたブロックリストは、複製によって更新されたブロックのリストを含む。

20

#### 【0062】

いくつかの実施の形態では、新たな情報源は、更新されたブロックリストに基づいて宛先へ送信されるデータを決定する(710)。情報源は、当該新たな情報源及び宛先を同期させるのに十分なデータを宛先にコピーする(712)。

30

#### 【0063】

一般的なインプリメンテーションでは、新たな情報源又は新たなハブは、新たな情報源を特定するコマンドを送信し、そのストレージレイが新たなハブとして引き継いでいることを示す。また、宛先間の前の相互通信が、ベースラインとなる1組のシーケンス番号を供給しているので、このコマンドは、前に特定されていなかった未処理のブロック書き込みのリストを特定する最終シーケンス番号のリストを送信するように各宛先に要求する。したがって、再形成のための相互通信は削減される。また、このコマンドは、宛先が旧情報源からのあらゆる新たな書き込みの受け取りを中止することも指定する。

#### 【0064】

図7Bでは、ファンアウト構成は、情報源の障害発生時の最新の宛先のステータスに再形成される。レイのステータスが求められる(714)。要求が、最新の状態を有する宛先へ送信される。再形成コマンド後、新たな情報源は、どのレイが最新であるかを判断するのに十分な情報を有する。最新のレイは、ローカルなブロックテーブルにおいて最も高いシーケンス番号を有するレイとして定義される。多くの場合、複数のレイが、レイの集合全体において「最新」である同一の状態を有する場合がある。それら複数のレイの1つを、最新として機能するように選択することができる。また、新たな情報源は、再形成のためにどのデータブロックを収集するかを判断するのに十分な情報も有する。新たな情報源は、最新の状態に達するのに十分なデータブロックを要求して(716)、フェッチし(718)、そして、要求したデータで当該新たな情報源を更新する(7

40

50

20)。したがって、情報源は、当該情報源に含まれるアレイの更新を制御する。

【0065】

図8A及び図8Bを参照して、概略ブロック図が、情報の追跡及び追跡された情報の共有を含まないストレージシステム配置800を示している。この例示の配置800は、中央のハブ802S、及び、ハブ802Sから1つ又は複数の宛先アレイ802Dへ放射状に伸びる通信スポークを備えた車輪として(as a wheel with)把握することができる。ハブ802Sは、情報源論理ユニット(lun)810Sを含むアレイとすることができる。スポークは、ハブ802Sを宛先アレイ802Dに接続する通信リンク808である。宛先アレイ802Dはリモートlun810Dを含む。データは、各通信リスト808上を同期又は非同期のいずれかで流れることができる。

10

【0066】

図8Aに示すように、ハブアレイ802Sが喪失する状態もしくは状況、又は、ハブ802Sへの通信が喪失する状態もしくは状況が発生すると、環境800は、オペレーションを宛先アレイ802Dの1つにフェイルオーバーして、アプリケーションの継続を可能にするように望まれる。従来のファンアウト技術では、宛先アレイ802D間の関連が存在しない。各宛先アレイ802Dは、ハブ802Sとの関係を有するだけである。

【0067】

図8Bに示すように、ファンアウト関係は、ハブ802Sの喪失のために、再形成を試みる。宛先アレイ802Dは、他の宛先アレイのlun810Dにどのブロックが書き込まれて、どのブロックが書き込まれていないのかに関する情報を含まない。その結果、新たなハブ802S'としてオペレーションを開始するように決定された宛先アレイは、lun810S'を他の宛先アレイ802Dのそれぞれに全面的にコピーしなければならず、これは非常に大きくなるおそれがある。

20

【0068】

再形成時に全面的なコピーを伴う技法と比較して、例示の構造及び技法では、複製の効率性が改善され、また、全面的なコピーを必要としない技法と比較しても、複製の効率性が改善される。

【0069】

例示の構造及び技法によって、任意の宛先を、新たな情報源として機能するように選択することが可能になる。

30

【0070】

本開示は、さまざまな実施の形態を説明しているが、これらの実施の形態は、例示として理解されるべきであって、特許請求の範囲を限定するものではない。説明した実施の形態の多くの変形、変更、追加、及び改良が可能である。たとえば、当業者は、本明細書で開示した構造及び方法を提供するのに必要なステップを容易に実施し、プロセスパラメータ、素材、及び寸法は単なる例としてのみ与えられていることを理解する。パラメータ、素材、及び寸法は、所望の構造及び変更を達成するように変化させることができ、これらの構造及び変更も特許請求の範囲内にある。本明細書で開示した実施の形態の変形及び変更も、添付の特許請求の範囲内に留まりつつ行うことができる。たとえば、開示した装置及び技法は、任意の適切な個数の記憶素子を有するストレージ構成において使用することができる。lunファンアウトは、例示として1:3ファンアウトとして示されている。例示の構造及び技法を使用すると、任意の適切なファンアウト比をサポートすることができる。ストレージシステムは、通常、磁気ディスク記憶素子を備えるが、任意の適切なタイプのストレージ技術を実施することができる。システムは、さまざまなオペレーティングシステム及びデータベースシステムで実施することができる。制御素子は、汎用のコンピュータシステム、ワークステーション、サーバ等でソフトウェア又はファームウェアとして実施することができるが、それ以外に、専用デバイス及び組み込みシステムにおいて実施することもできる。

40

【図面の簡単な説明】

【0071】

50

【図1A】1:nファンアウト構成で配置されて、少ないオーバーヘッドで1:n-1ファンアウトに修復するのに適合した、複数のストレージアレイを有するストレージシステムの一実施の形態を示す概略ブロック図である。

【図1B】1:nファンアウト構成で配置されて、少ないオーバーヘッドで1:n-1ファンアウトに修復するのに適合した、複数のストレージアレイを有するストレージシステムの一実施の形態を示す概略ブロック図である。

【図1C】1:nファンアウト構成で配置されて、少ないオーバーヘッドで1:n-1ファンアウトに修復するように適合された、複数のストレージアレイを有するストレージシステムの一実施の形態を示す概略ブロック図である。

【図2】冗長データストレージシステムでの使用に適合されたストレージユニットの一実施の形態を示す概略ブロック図である。 10

【図3】危険性なくデータを維持しつつそのデータを複数の地理的位置に複製するように、ファンアウト関係を高速に再形成するよう適合された技法の一実施の形態を示す概略フローチャートである。

【図4】ストレージシステムのデータ追跡の一例におけるブロックマップシーケンスを示す概略図である。

【図5】複製追跡中にデータを収集するための使用に適したデータ構造体の別の実施の形態を示す概略テーブル図である。

【図6A】情報源の障害の場合に、ファンアウトを削減されたファンアウト比に修復するための技法の実施の形態を示すフローチャートである。 20

【図6B】情報源の障害の場合に、ファンアウトを削減されたファンアウト比に修復するための技法の実施の形態を示すフローチャートである。

【図7A】情報源の障害の発生時にファンアウト構成を再形成するための技法の実施の形態を示すフローチャートである。

【図7B】情報源の障害の発生時にファンアウト構成を再形成するための技法の実施の形態を示すフローチャートである。

【図8A】情報の追跡及び追跡された情報の共有を含まないストレージシステムの配置を示す概略ブロック図である。

【図8B】情報の追跡及び追跡された情報の共有を含まないストレージシステムの配置を示す概略ブロック図である。 30

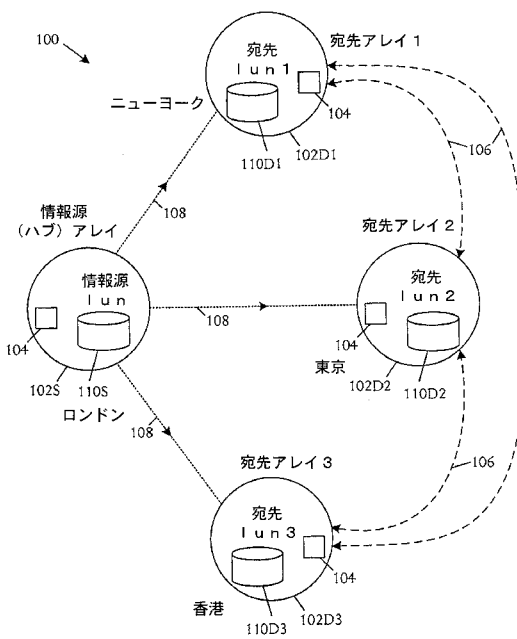
#### 【符号の説明】

#### 【0072】

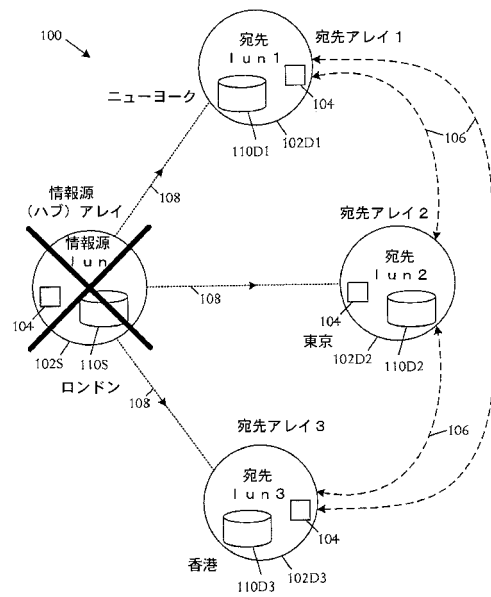
100・・・ストレージシステム  
 102S・・・情報源(ハブ)アレイ  
 102D1・・・宛先アレイ1  
 102D2・・・宛先アレイ2  
 102D3・・・宛先アレイ3  
 110S・・・情報源1un  
 110D1・・・宛先1un1  
 110D2・・・宛先1un2  
 110D3・・・宛先1un3  
 102S'・・・新たな情報源(ハブ)アレイ  
 102D1'・・・宛先アレイ1'  
 102D2'・・・宛先アレイ2'  
 110S'・・・新たな情報源1un  
 110D1'・・・宛先1un1'  
 110D2'・・・宛先1un2'  
 202・・・ストレージユニット  
 204・・・ストレージ  
 206・・・I/Oインターフェイス 40

- 208・・・コントローラ
- 802S・・・情報源（ハブ）アレイ
- 802D・・・宛先アレイ1、2、3
- 810S・・・情報源 lun
- 810D・・・宛先 lun 1、2、3
- 802S'・・・新たな情報源（ハブ）アレイ
- 802D'・・・宛先アレイ1'
- 802D'・・・宛先アレイ2'
- 810S'・・・情報源 lun
- 810S'・・・新たな情報源 lun
- 810D'・・・宛先 lun 1'
- 810D'・・・宛先 lun 2'

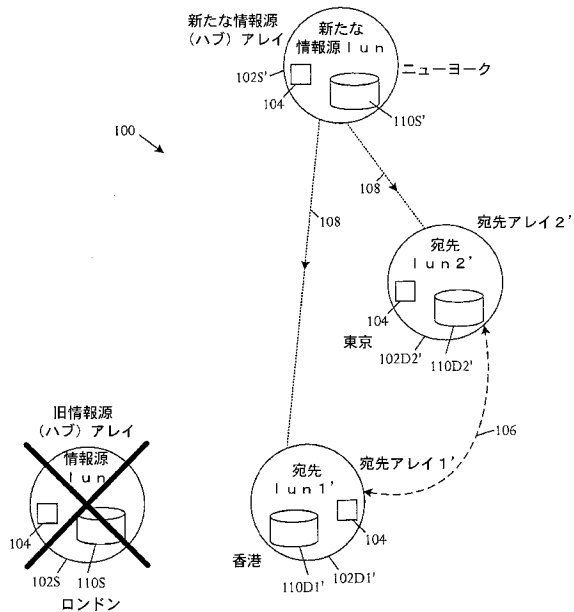
【図1A】



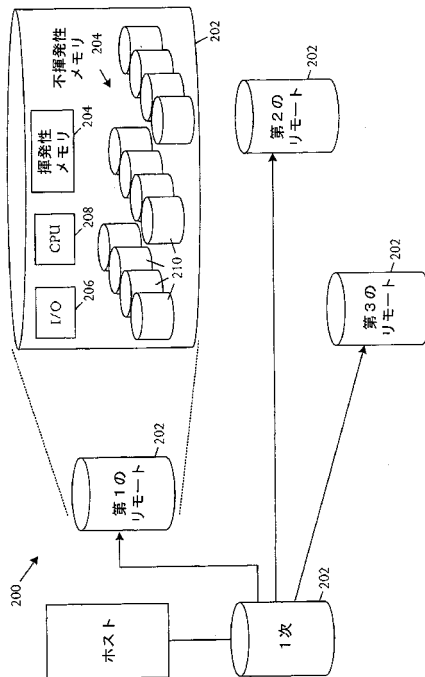
【図1B】



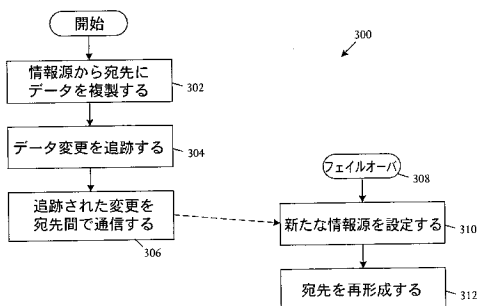
【図1C】



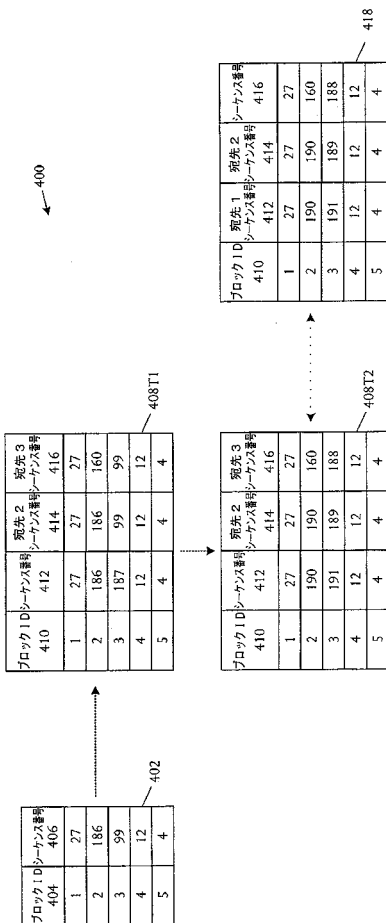
【図2】



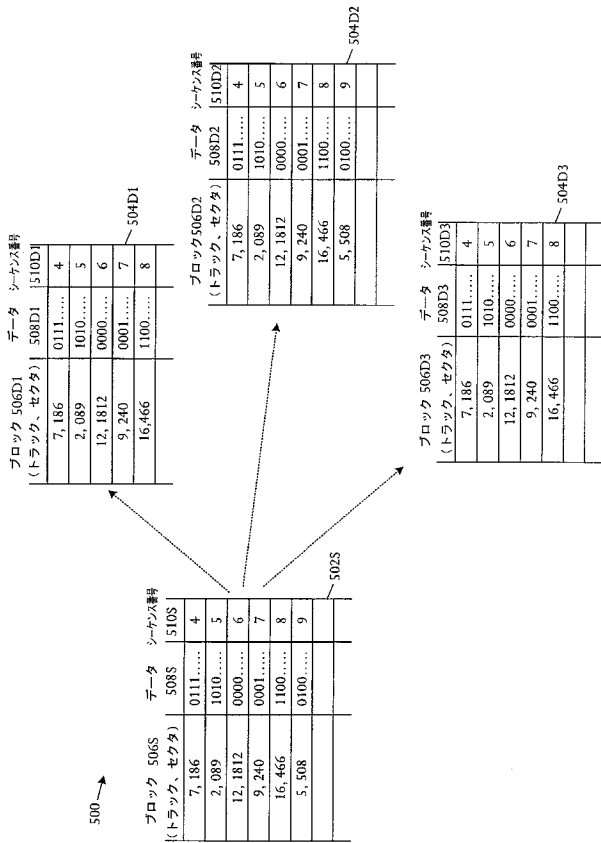
【図3】



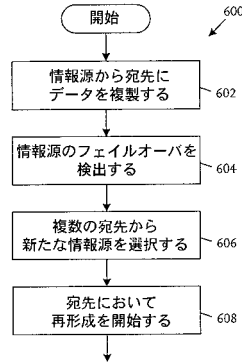
【図4】



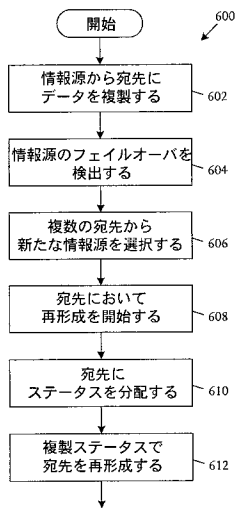
【 図 5 】



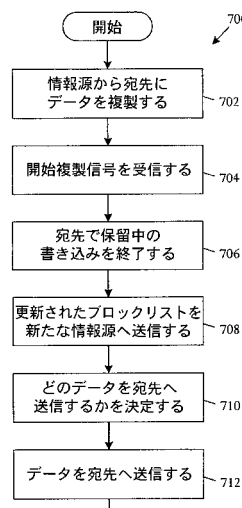
【 図 6 A 】



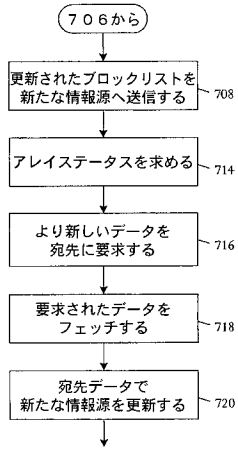
【 図 6 B 】



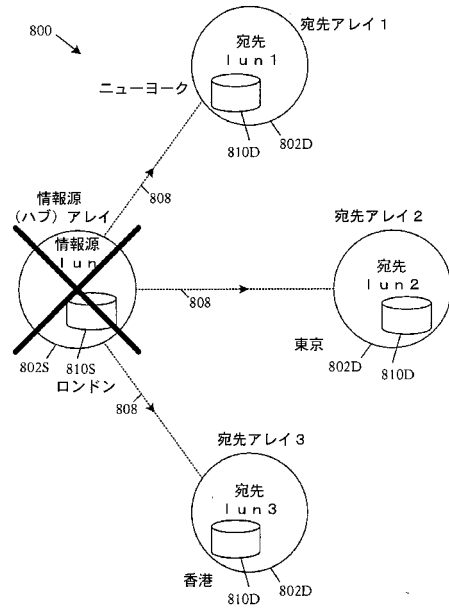
【 図 7 A 】



【 図 7 B 】



【 図 8 A 】



【 図 8 B 】

