

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2017-531250

(P2017-531250A)

(43) 公表日 平成29年10月19日(2017.10.19)

(51) Int.Cl.  
G06F 12/00 (2006.01)

F I  
G06F 12/00 533 J

テーマコード (参考)

審査請求 有 予備審査請求 未請求 (全 42 頁)

(21) 出願番号 特願2017-511735 (P2017-511735)  
 (86) (22) 出願日 平成27年7月31日 (2015.7.31)  
 (85) 翻訳文提出日 平成29年4月14日 (2017.4.14)  
 (86) 国際出願番号 PCT/US2015/043159  
 (87) 国際公開番号 WO2016/032688  
 (87) 国際公開日 平成28年3月3日 (2016.3.3)  
 (31) 優先権主張番号 14/473, 621  
 (32) 優先日 平成26年8月29日 (2014.8.29)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 303039534  
 ネットアップ, インコーポレイテッド  
 アメリカ合衆国 カリフォルニア 940  
 89, サニーヴェール, イースト ジャ  
 ヴァ ドライブ 495  
 (74) 代理人 100107766  
 弁理士 伊東 忠重  
 (74) 代理人 100070150  
 弁理士 伊東 忠彦  
 (74) 代理人 100091214  
 弁理士 大貫 進介  
 (72) 発明者 キンメル, ジェフリー, エス.  
 アメリカ合衆国 カリフォルニア州 94  
 089 サニーベール イースト ジャ  
 ヴァ ドライブ 495

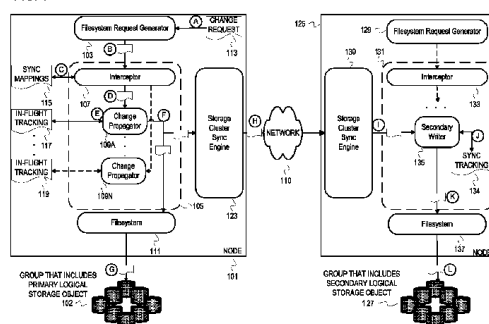
最終頁に続く

(54) 【発明の名称】 粒状同期/半同期アーキテクチャ

(57) 【要約】

データ一貫性及び可用性が、クラスタ化された記憶環境内で記憶仮想化を使用する記憶ソリューションにおいて、論理記憶オブジェクトの粒度で提供されることができる。異なる記憶要素にわたりデータの一貫性を確保するために、異なる記憶要素にわたって同期が実行される。データに対する変更が、該変更をプライマリ論理記憶要素からセカンダリ論理記憶要素に伝搬することによって、異なるクラスタ内の記憶要素にわたり同期される。パフォーマンスを維持すると同時に最も厳しいRPOを満足するために、変更要求が、プライマリ論理記憶オブジェクトをホストするファイルシステムに送信される前に傍受され、セカンダリ論理記憶オブジェクトに関連付けられた異なる管理記憶要素に伝搬される。

FIG. 1



**【特許請求の範囲】****【請求項 1】**

記憶プロトコル入力/出力要求からファイルシステム要求を生成した後、前記記憶プロトコル入力/出力要求の対象として示される第 1 の論理記憶オブジェクトが同期構成データ内に示されるとおりに第 2 の論理記憶オブジェクトとの同期関係を有すると決定することと、

前記同期構成データ内の前記同期関係について、複数の同期関係タイプのうちいずれが示されているかを決定することと、

前記複数の同期関係タイプのうちの第 1 の同期関係タイプについて、

前記ファイルシステム要求を追跡するための前記ファイルシステム要求の指標を記録することと、

前記ファイルシステム要求をファイルシステムに供給することであって、前記ファイルシステムは、前記第 1 の論理記憶オブジェクトをホストする 1 つ以上の記憶要素にアクセスする、ことと、

前記ファイルシステム要求に基づいて要求を生成することであって、前記の生成された要求は前記第 2 の論理記憶オブジェクトを対象として示す、ことと、

前記の生成された要求を伝達のために前記第 2 の論理記憶オブジェクトに関連付けられたノードに供給することと、

前記記憶プロトコル入力/出力要求により示される変更が前記第 1 の論理記憶オブジェクトに対して成功裏に行われ、前記変更が前記第 2 の論理記憶オブジェクトに対して成功裏に行われる場合、前記記憶プロトコル入力/出力要求の成功を示す応答を前記記憶プロトコル入力/出力要求の要求元に供給することと、

を含む方法。

**【請求項 2】**

前記複数の同期関係タイプのうちの第 2 の同期関係タイプについて、

所定の時間間隔において、前記第 1 の論理記憶オブジェクトを対象にする複数のファイルシステム要求を蓄積することであって、前記複数のファイルシステム要求は前記ファイルシステム要求を含む、ことと、

前記複数のファイルシステム要求の各々を前記ファイルシステムに供給することと、

前記の蓄積された複数のファイルシステム要求のうち、前記第 1 の論理記憶オブジェクトに対して行われる変更を示すファイルシステム要求間の依存関係を決定することと、

シーケンシング情報を記録して、決定された依存関係を保存することと、

前記第 2 の論理記憶オブジェクトに関連付けられた前記ノードに、前記シーケンシング情報と、変更を示す前記の蓄積された複数のファイルシステム要求とを供給することであって、該変更は前記第 2 の論理記憶オブジェクトに対して行われる、ことと、

をさらに含む請求項 1 に記載の方法。

**【請求項 3】**

前記第 1 の論理記憶オブジェクトに対して行われる変更を示す前記複数のファイルシステム要求のうち冗長なファイルシステム要求を消去すること、をさらに含む請求項 2 に記載の方法。

**【請求項 4】**

前記第 1 の論理記憶オブジェクトに対して行われる変更を示しかつ前記所定の時間間隔内に蓄積された前記複数のファイルシステム要求の総数を示すメタデータを供給すること、をさらに含む請求項 2 又は 3 に記載の方法。

**【請求項 5】**

前記第 2 の論理記憶オブジェクトに関連付けられた前記ノードに供給される前記複数のファイルシステム要求の各々に前記メタデータを記録すること、をさらに含む請求項 4 に記載の方法。

**【請求項 6】**

前記シーケンシング情報を記録することは、依存関係を有しかつ前記第 1 の論理記憶オ

10

20

30

40

50

プロジェクトに対して行われる変更を示す前記複数のファイルシステム要求の各々に前記シーケンシング情報を記録することを含む、請求項 2 乃至 5 のうちいずれか 1 項に記載の方法。

【請求項 7】

前記記憶プロトコル入力/出力要求により示される変更が前記第 1 の論理記憶オブジェクト又は前記第 2 の論理記憶オブジェクトのいずれかに対して成功裏に行われない場合、前記記憶プロトコル入力/出力要求の失敗を示す応答を前記記憶プロトコル入力/出力要求の前記要求元に供給すること、をさらに含む請求項 1 乃至 6 のうちいずれか 1 項に記載の方法。

【請求項 8】

1 つ以上のコンピュータシステムにより実行されるときに、前記 1 つ以上のコンピュータシステムに請求項 1 乃至 7 のうちいずれか 1 項に記載の方法を実施させるマシン読取可能命令を含むコンピュータプログラム。

【請求項 9】

記憶クラスタにわたる論理記憶オブジェクト間の同期関係を維持するコンピュータプログラムであって、当該コンピュータプログラムは、非一時的マシン読取可能媒体内に具現化されるプログラムコードを含み、前記プログラムコードは、

記憶プロトコル入力/出力要求からのファイルシステム要求が生成された後、前記記憶プロトコル入力/出力要求の対象として示される第 1 の論理記憶オブジェクトが同期構成データ内に示されるとおりに第 2 の論理記憶オブジェクトとの同期関係を有すると決定し

、前記同期構成データ内の前記同期関係について、複数の同期関係タイプのうちいずれが示されているかを決定し、

前記複数の同期関係タイプのうちの第 1 の同期関係タイプについて、

前記ファイルシステム要求を追跡するための前記ファイルシステム要求の指標を記録し、

前記ファイルシステム要求をファイルシステムに供給することであって、前記ファイルシステムは、前記第 1 の論理記憶オブジェクトをホストする 1 つ以上の記憶要素にアクセスし、

前記ファイルシステム要求に基づいて要求を生成することであって、前記の生成された要求は前記第 2 の論理記憶オブジェクトを対象として示し、

前記の生成された要求を伝達のために前記第 2 の論理記憶オブジェクトに関連付けられたノードに供給し、

前記記憶プロトコル入力/出力要求により示される変更が前記第 1 の論理記憶オブジェクトに対して成功裏に行われ、前記変更が前記第 2 の論理記憶オブジェクトに対して成功裏に行われる場合、前記記憶プロトコル入力/出力要求の成功を示す応答を前記記憶プロトコル入力/出力要求の要求元に供給する

ように構成される、コンピュータプログラム。

【請求項 10】

前記複数の同期関係タイプのうちの第 2 の同期関係タイプについて、

所定の時間間隔において、前記第 1 の論理記憶オブジェクトを対象にする複数のファイルシステム要求を蓄積することであって、前記複数のファイルシステム要求は前記ファイルシステム要求を含み、

前記複数のファイルシステム要求の各々を前記ファイルシステムに供給し、

前記の蓄積された複数のファイルシステム要求のうち、前記第 1 の論理記憶オブジェクトに対して行われる変更を示すファイルシステム要求間の依存関係を決定し、

シーケンシング情報を記録して、決定された依存関係を保存し、

前記第 2 の論理記憶オブジェクトに関連付けられた前記ノードに、前記シーケンシング情報と、変更を示す前記の蓄積された複数のファイルシステム要求とを供給することであって、該変更は前記第 2 の論理記憶オブジェクトに対して行われる

10

20

30

40

50

ように構成されたプログラムコード、をさらに含む請求項 9 に記載のコンピュータプログラム。

【請求項 1 1】

前記第 1 の論理記憶オブジェクトに対して行われる変更を示す前記複数のファイルシステム要求のうち冗長なファイルシステム要求を消去するプログラムコード、をさらに含む請求項 1 0 に記載のコンピュータプログラム。

【請求項 1 2】

前記第 1 の論理記憶オブジェクトに対して行われる変更を示しかつ前記所定の時間間隔内に蓄積された前記複数のファイルシステム要求の総数を示すメタデータを供給するプログラムコード、をさらに含む請求項 1 0 又は 1 1 に記載のコンピュータプログラム。

10

【請求項 1 3】

前記第 2 の論理記憶オブジェクトに関連付けられた前記ノードに供給されることになる前記複数のファイルシステム要求の各々に前記メタデータを記録するプログラムコード、をさらに含む請求項 1 2 に記載のコンピュータプログラム。

【請求項 1 4】

前記シーケンシング情報を記録する前記プログラムコードは、依存関係を有しかつ前記第 1 の論理記憶オブジェクトに対して行われる変更を示す前記複数のファイルシステム要求の各々に前記シーケンシング情報を記録するプログラムコードを含む、請求項 1 0 乃至 1 3 のうちいずれか 1 項に記載のコンピュータプログラム。

【請求項 1 5】

20

前記記憶プロトコル入力/出力要求により示される変更が前記第 1 の論理記憶オブジェクト又は前記第 2 の論理記憶オブジェクトのいずれかに対して成功裏に行われない場合、前記記憶プロトコル入力/出力要求の失敗を示す応答を前記記憶プロトコル入力/出力要求の前記要求元に供給するプログラムコード、をさらに含む請求項 9 乃至 1 4 のうちいずれか 1 項に記載のコンピュータプログラム。

【請求項 1 6】

装置であって、  
プロセッサと、  
プログラムコードを記憶させた非一時的マシン読取可能媒体と、  
を含み、

30

前記プログラムコードは前記プロセッサにより実行可能であって、当該装置に、

記憶プロトコル入力/出力要求からのファイルシステム要求が生成された後、前記記憶プロトコル入力/出力要求の対象として示される第 1 の論理記憶オブジェクトが同期構成データ内に示されるとおりに第 2 の論理記憶オブジェクトとの同期関係を有すると決定し、

前記同期構成データ内の前記同期関係について、複数の同期関係タイプのうちいずれが示されているかを決定し、

前記複数の同期関係タイプのうちの第 1 の同期関係タイプについて、

前記ファイルシステム要求を追跡するための前記ファイルシステム要求の指標を記録し、

40

前記ファイルシステム要求をファイルシステムに供給することであって、前記ファイルシステムは、前記第 1 の論理記憶オブジェクトをホストする 1 つ以上の記憶要素にアクセスし、

前記ファイルシステム要求に基づいて要求を生成することであって、前記の生成された要求は前記第 2 の論理記憶オブジェクトを対象として示し、

前記の生成された要求を伝達のために前記第 2 の論理記憶オブジェクトに関連付けられたノードに供給し、

前記記憶プロトコル入力/出力要求により示される変更が前記第 1 の論理記憶オブジェクトに対して成功裏に行われ、前記変更が前記第 2 の論理記憶オブジェクトに対して成功裏に行われる場合、前記記憶プロトコル入力/出力要求の成功を示す応答を前記記憶プ

50

ロトコル入力/出力要求の要求元に供給することを  
実行させる、装置。

【請求項 17】

前記プログラムコードは、当該装置に、  
前記複数の同期関係タイプのうちの第2の同期関係タイプについて、  
所定の時間間隔において、前記第1の論理記憶オブジェクトを対象にする複数のファイルシステム要求を蓄積することであって、前記複数のファイルシステム要求は前記ファイルシステム要求を含み、  
前記複数のファイルシステム要求の各々を前記ファイルシステムに供給し、  
前記の蓄積された複数のファイルシステム要求のうち、前記第1の論理記憶オブジェクトに対して行われる変更を示すファイルシステム要求間の依存関係を決定し、  
シーケンシング情報を記録して、決定された依存関係を保存し、  
前記第2の論理記憶オブジェクトに関連付けられた前記ノードに、前記シーケンシング情報と、変更を示す前記の蓄積された複数のファイルシステム要求とを供給することであって、該変更は前記第2の論理記憶オブジェクトに対して行われる  
ことを実行させるための、前記プロセッサにより実行可能なプログラムコードをさらに含む、請求項16に記載の装置。

10

【請求項 18】

前記プログラムコードは、前記第1の論理記憶オブジェクトに対して行われる変更を示す前記複数のファイルシステム要求のうち冗長なファイルシステム要求を消去することを当該装置に実行させるための、前記プロセッサにより実行可能なプログラムコードをさらに含む、請求項17に記載の装置。

20

【請求項 19】

前記プログラムコードは、前記第1の論理記憶オブジェクトに対して行われる変更を示しかつ前記所定の時間間隔内に蓄積された前記複数のファイルシステム要求の総数を示すメタデータを供給することを当該装置に実行させるための、前記プロセッサにより実行可能なプログラムコードをさらに含む、請求項17又は18に記載の装置。

【請求項 20】

前記プログラムコードは、前記記憶プロトコル入力/出力要求により示される変更が前記第1の論理記憶オブジェクト又は前記第2の論理記憶オブジェクトのいずれかに対して成功裏に行われない場合、前記記憶プロトコル入力/出力要求の失敗を示す応答を前記記憶プロトコル入力/出力要求の前記要求元に供給することを当該装置に実行させるための、前記プロセッサにより実行可能なプログラムコードをさらに含む、請求項16乃至19のうちいずれか1項に記載の装置。

30

【発明の詳細な説明】

【技術分野】

【0001】

関連出願

本出願は、“GRANULAR SYNC/SEMI-SYNC ARCHITECTURE”と題され2014年8月29日に申請された米国特許出願第14/473,621号に対する優先を主張し、上記出願は本明細書において参照により援用される。

40

【0002】

本開示の態様は、概して分散記憶の分野に関し、より詳細には分散ストレージにわたりデータを同期するアーキテクチャに関する。

【背景技術】

【0003】

企業は、顧客データを維持するか又は企業独自データを維持するかにかかわらず、常に利用可能な又は高度に利用可能なデータと、そのデータの保護とを求める。こうした需要をサポートするために、データはしばしば、複数のサイトにおける複数の記憶システムに

50

わたり存在し、複数のサイトはしばしば、大きく距離が離れている。上記サイトが長距離離れている理由の1つは、1つの大惨事がデータ可用性に影響するのを回避するためである。可用性要件を定義するのに使用されるメトリクスには、復旧ポイント目標（RPO）及び復旧時間目標（RTO）が含まれる。企業は、RTOを、企業が企業データへのアクセスの欠如を許容する最大量の時間として指定する。企業は、RPOを、中断に起因して失われる可能性のある、時間の観点からのデータの量として指定する。例えば、企業は、RTOを15秒として指定することができる。換言すると、企業は、サービス中断又は失敗の時間からそのシステムの完全復旧の時間まで、最大で15秒を受け入れることになる。RPOについて、企業は、5秒を指定することができる。このことは、企業が、失敗又は中断の前の5秒内に書き込まれたデータ（例えば、新しい書き込み、更新等）を超える

10

#### 【0004】

記憶システムにわたる企業の可用性及び保護の需要をサポートするための記憶機構は、様々な名称を付与されており、例えば、スナップショット（snapshotting）、ミラーリング、クローニング、及びレプリケーション（replicating）などである。上記記憶機構の各々は、ストレージ機構及び/又はストレージ製品のプロバイダによってさらに変動することがある。こうした変動にかかわらず、各記憶機構は、企業のデータについての一貫性のあるビューを提供する。

#### 【図面の簡単な説明】

#### 【0005】

20

本開示の態様は、添付図面を参照することによって、より理解され得る。

【図1】変更要求に応答して異なるクラスタ内のプライマリ論理記憶オブジェクトとセカンダリ論理記憶オブジェクトとの間のデータ交換を調整する例示的な記憶クラスタ同期エンジンを表す。

【図2】変更要求に応答して異なるクラスタ内のプライマリ論理記憶オブジェクトとセカンダリ論理記憶オブジェクトとの間のデータ交換を調整する例示的な記憶クラスタ同期エンジンを表す。

【図3】完全同期関係においてセカンダリ論理記憶オブジェクトとして構成され、及び半同期関係においてプライマリ論理記憶オブジェクトとして構成された論理記憶オブジェクトの、例示的な動作を表す。

30

【図4】完全同期関係においてセカンダリ論理記憶オブジェクトとして構成され、及び半同期関係においてプライマリ論理記憶オブジェクトとして構成された論理記憶オブジェクトの、例示的な動作を表す。

【図5】1つ又は複数のクラスタのノードにわたり論理記憶オブジェクト粒度において完全同期及び半同期の双方を提供する一例示的なアーキテクチャを表す。

【図6】変更要求の受信を扱うことと完全同期関係におけるプライマリエンドポイントを対象にする変更要求を扱うこととの例示的な動作のフローチャートを表す。

【図7】論理記憶オブジェクト粒度半同期動作についての例示的な動作のフローチャートを表す。

【図8】クローズされた変更設定ログを処理する例示的な動作のフローチャートを表す。

40

【図9】同期関係におけるプライマリエンドポイントの記憶要素モジュールからの応答を扱う例示的な動作のフローチャートを表す。

【図10】クラスタベースの同期エンジンが伝搬器とカウンターパートの同期エンジンからの要求を処理する例示的な動作のフローチャートを表す。

【図11】伝搬器インスタンスがセカンダリエンドポイントに対する変更要求への応答を扱う例示的な動作のフローチャートを表す。

【図12】セカンダリライタが複製要求を扱う例示的な動作のフローチャートを表す。

【図13】セカンダリライタインスタンスが下層の記憶要素モジュールからの応答を扱う例示的な動作のフローチャートを表す。

【図14】記憶クラスタベース粒状完全同期及び半同期伝搬エンジンを備えた一例示的な

50

コンピュータシステムを表す。

【発明を実施するための形態】

【0006】

下記の説明は、本開示の手法を具現化する例示的なシステム、方法、手法、命令シーケンス、及びコンピュータプログラム製品を含む。しかしながら、説明される本開示の態様は、こうした特定の詳細なしに実施され得ることが理解される。例えば、例示がディスク及びディスクアレイを参照するが、開示の態様はそのように限定されない。開示の態様は、ソリッドステート記憶装置、光学記憶装置、個々の記憶装置の連合体、異なるタイプの記憶装置の組み合わせ等を使用する記憶システムに実装され得る。さらに、多くの例示が、論理記憶オブジェクトのペアを使用して動作を例示する。開示の態様は、プライマリ（primary）及びセカンダリ（secondary）論理記憶オブジェクトペアに限定されず、論理記憶オブジェクトのグループに適用されることができる。例えば、プライマリ論理記憶オブジェクトを複数のセカンダリ論理記憶オブジェクトと同期するようにシステムが構成されることができる。良く知られる命令インスタンス、プロトコル、構造、及び手法は、説明を分かりにくくしないように詳細には図示されていない。

10

【0007】

#### 用語

本説明は、用語「記憶要素」を使用して、データをホストし及び/又はデータに対するアクセスを管理する記憶システム内の任意のエンティティを参照する。本明細書において参照される記憶要素は、管理記憶要素及びホスト記憶要素として分類されることができる。管理記憶要素とホスト記憶要素との間の区別は、記憶要素の基本機能性（primary functionality）から生じる。管理記憶要素は、ホスト記憶要素に対するアクセスを主に管理する。管理記憶要素は、他の装置（例えば、クライアント）からの要求を処理し、動作（例えば、スナップショット動作）を実行するための要求を発する（originate）ことができる。要求が別の装置からであるか又は管理記憶要素から発せられるかにかかわらず、管理記憶要素は、要求をホスト記憶要素に伝達する。管理記憶要素の例には、ファイルサーバ及び記憶コントローラが含まれる。ホスト記憶要素は、管理記憶要素の観点から要求を最終的に満たす動作を主に実行する。ホスト記憶要素は、管理記憶要素からの要求により指定される場所の読み出し又は該場所への書き込みを実行する。この読み出し又は書き込みは、一ディスク又は複数ディスクに対して実行されてよい。仮想化の複数レイヤの場合、読み出し又は書き込みは、管理記憶要素の観点から1つ又は複数のディスクであるように見えるものに対して実行されてよい。ホスト記憶要素の例には、ディスクドライブ、光学ドライブ、記憶アレイ、及びテープドライブが含まれる。

20

30

【0008】

用語の管理記憶要素及びホスト記憶用は記憶要素の基本機能性に基づいて使用され、なぜならば、機能性が要素間で排他的ではないからである。例えば、記憶コントローラが、キャッシュ内にローカルに記憶されたデータを有して、アクセス要求の扱いを促進してもよい。記憶コントローラがアクセス要求を満たすことができるとしても、記憶コントローラの基本機能性は、ローカルメモリからデータを読み出すことやローカルメモリにデータを書き込むことではない。同様に、ホスト記憶要素が、ディスクに対するアクセスを管理するハードウェアを含むことができる。例えば、RAID（redundant array of independent disks）コントローラとディスクのアレイとが、単一の筐体内に収納されることができる。RAIDコントローラはディスクのアレイに対するアクセスを管理するが、上記単一筐体内に収納されるコンポーネントの基本機能性は、管理記憶要素から受信される要求を満たすことである。

40

【0009】

本説明は、用語の完全同期（full synchronization）（「完全同期（full sync）」）及び半同期（semi-synchronization）（「半同期（semi sync）」）をさらに使用する。上記用語は、異なるタイプの同期構成を参照する。「完全同期」構成は、本明細書において使用されるとき、変更要求が実行されたことを確認するリプライの送信を、該変更が

50

プライマリ論理記憶オブジェクトとセカンダリ論理記憶オブジェクトとにわたり同期されるまで遅延させる構成を参照する。「半同期」構成は、本明細書において使用されるとき、変更要求が実行されたことを確認するリプライが、プライマリ論理記憶オブジェクトに対して該変更が実行された後、(1つ以上の)セカンダリ論理記憶オブジェクトとの同期が依然として進行中であり得る間に送信されることを可能にする構成を参照する。

#### 【0010】

本説明は、用語「要求」を使用して、何らかが行われることを要求するソフトウェアエンティティ又はハードウェアエンティティ間における通信を参照し、様々なプロトコルにおいて使用される名称、データフィールド等のパリエーションを回避する。要求は、データが読み出されること、データが書き込まれることの要求、又は何らかの他のデータ処理要求を示すことができる。要求は、動作のタイプ(例えば、読み出し、書き込み)、要求の対象(例えば、論理記憶オブジェクト識別子)、及び要求元(requestor)の識別子を示すことができる。さらなる情報が、統制するプロトコルに依存して要求の中に示されてもよい。しかし、本説明は、こうしたさらなる情報の詳細を掘り下げない。さらに、複数のプロトコルが、プロトコルスタックといわれるものを形成することができる。プロトコルスタックは、要求が通過し又は横断する一連の処理モジュールとみなされることができる。プロトコルスタックの各レイヤにおいて、ヘッダ及び/又はトレーラが要求に追加され又は要求から削除されてもよい。本説明において、少なくともいくつかのスタック処理は、説明にさらなる複雑さを加えるのを回避するために、説明されない。本説明は、要求を、関連付けられたヘッダ又はトレーラにかかわらない、かつヘッダ及び/又はトレーラ内の値に対するとり得る修正にかかわらない要求として、参照することになる。

10

20

#### 【0011】

##### 導入

クラスタリングは、一般に、ハードウェア要素と一緒にグループ化することを参照し、個々のハードウェア要素から得ることができないハードウェア要素(例えば、ディスクドライブ、記憶アレイ、ファイルサーバ、記憶コントローラ等)のグループ(「クラスタ」)の恩恵を得る。クラスタリングは、様々な記憶機構に対して使用することができ、その例には、ロードバランシング、フェイルオーバーサポート、I/O帯域幅の増大、及びデータ可用性が含まれる。こうした記憶の態様をサポートし、一貫性のある記憶のビューを提供するために、サポート記憶要素間でデータが同期される。データを要求元に供給し、要求元のためにデータを修正するのに、いずれの記憶要素が(例えば、構成によって)最初に及び/又は優先的に使用されるかに依存して、種々のハードウェア記憶要素がプライマリ記憶要素及びセカンダリ記憶要素としばしばいわれる。さらに、記憶要素のークラスタがプライマリクラスタとして指名されることができ、記憶要素のークラスタがセカンダリクラスタとして指名されることができる。

30

40

#### 【0012】

多くの記憶システム機能性が、記憶仮想化の特徴として展開される。しばしば、記憶仮想化ソフトウェア/ツールが、記憶システムを構成する実際のハードウェア要素を分かりにくくする。ゆえに、要求元(ときに、本明細書においてクライアントといわれる)は、しばしば、論理記憶オブジェクト又は論理記憶コンテナから読み出し及び書き込みをし、その例には、論理ユニット番号(LUN)、ファイル、仮想マシンディスク(VMDK)、仮想ボリューム、及び論理パーティションが含まれる。仮想化の任意数のレイヤが、実際の記憶システムハードウェア要素と、アクセス要求を送信するクライアントとを分離することができる。各記憶システムハードウェア要素は、多数の論理記憶オブジェクト及び/又は論理記憶オブジェクトの多数の部分をホストすることができる。さらに、クライアントのために要求を扱う記憶コントローラが、物理記憶アレイであるように見える仮想記憶アレイと通信してもよい。ゆえに、記憶アレイであるかのように提示される論理記憶オブジェクトが、複数の論理記憶オブジェクトをホストするものとして提示されることができる。

#### 【0013】

50



## 概説

データ一貫性及び可用性が、クラスタ化された記憶環境内で記憶仮想化を使用する記憶ソリューションにおいて、論理記憶オブジェクトの粒度で提供されることができる。可用性について、データは、前述されたとおり、異なるサイトにおける異なる記憶要素上で維持される。異なる記憶要素にわたりデータの一貫性を確保するために、異なる記憶要素にわたって同期が実行される。論理記憶オブジェクトの粒度において、離れたサイトにおける異なる記憶要素にわたりデータが効率的に同期されることができ、なぜならば、少なくとも部分的に、同期されるデータの量がより小さくなり、データを運ぶネットワーク内の否定的インシデントの影響をより受けにくくなるからである。データに対する変更が、該変更をプライマリ論理記憶オブジェクト（すなわち、変更要求の中に指定された論理記憶オブジェクト）に関連付けられたノードからセカンダリ論理記憶オブジェクト（すなわち、同期についてプライマリ論理記憶オブジェクトに関連付けられた論理記憶オブジェクト）に伝搬することによって、異なるクラスタ内の記憶要素にわたり同期される。パフォーマンスを維持すると同時に最も厳しいRPO（例えば、RPO = 0）及びRTOを満足するために、ファイルシステム要求が、プライマリ論理記憶オブジェクトをホストするファイルシステム（「プライマリファイルシステム」）に送信される前に傍受され（intercepted）、セカンダリ論理記憶オブジェクトに関連付けられたノードのファイルシステム（「セカンダリファイルシステム」）に伝搬される。論理記憶オブジェクトは、少なくともいくつかの関連クラスタ内で排他的である不変識別子を有して、クラスタにわたり論理記憶オブジェクトの有効な識別を可能にする。ファイルシステム要求が記憶プロトコル固有要求から生成された後、それがプライマリファイルシステムに送信される前に該ファイルシステム要求を傍受することは、記憶プロトコル固有及び/又はアプリケーション固有の動作を伝搬動作に負わせることを回避し、このことは、伝搬される変更要求のサイズと処理動作の数とをさらに低減させる。ファイルシステムと直接インターフェースをとる同期をサポートする動作を扱うエンティティを有することで、ファイルシステム応答の効率的な伝達についてファイルシステムのメカニズムを活用する。

10

20

30

40

50

### 【0014】

#### 例示

図1～4に表される例示は、システムに関して圧倒的な量の情報を提示するのを回避する試みとして、異なる度合の例示的な詳細を表す。あらゆるとり得るデータ構造とあらゆるとり得る機能性のモジュール化とが提示されるわけではなく、なぜならば、それらは多数であり、本開示の態様の理解に必要ではないからである。例えば、複数のデータ構造として提示されるデータ構造が、様々なとり得る索引/アクセススキームとデータの配置とを用いて、別様に編成されることができる。同様に、例示の中の個々のモジュール/エンジン/ユニットとして提示される機能性もまた、プラットフォーム（オペレーティングシステム及び/又はハードウェア）、アプリケーションエコシステム、インターフェース、プログラムの選好、プログラミング言語等のうち任意のものに従って、別様に編成されることができる。さらに、いくつかの機能性が、同様に圧倒的な量の情報を提示するのを回避する試みとして、本説明の中で後から説明される。例えば、管理エンティティからのスナップショット要求、又は半同期構成が、プライマリ管理記憶要素における複数のアクセス要求のシーケンシング（sequencing）につながることもある。シーケンシングは、前の方の例示の中では論じられない。

### 【0015】

図1及び図2は、変更要求に応答して異なるクラスタ内のプライマリ論理記憶オブジェクトとセカンダリ論理記憶オブジェクトとの間のデータ変更を調整する例示的な記憶クラスタ同期エンジンを表す。図1は、プライマリ管理記憶要素からセカンダリ管理記憶要素への変更要求の伝搬を表す。管理記憶要素は、簡潔さのため、以降はノードといわれる。図1において、第1のクラスタがプライマリノード101を含み、第2のクラスタがセカンダリノード125を含む。図の簡素化のため、及び描画空間制約に起因して、クラスタ全体は表されていない。プライマリノード101は、論理記憶オブジェクトのグループ1

02をホストするホスト記憶要素に通信可能に結合される。グループ102は、プライマリ論理記憶オブジェクトを含む。プライマリノード101は、ファイルシステム要求生成器103、変更伝搬エンジン105、記憶クラスタ同期エンジン123、及びファイルシステム111を含む。ファイルシステム要求生成器103は、記憶プロトコルベースの要求からファイルシステム要求を生成する。変更伝搬エンジン105は、傍受器(interceptor)107、変更伝搬器(change propagator)109A、及び変更伝搬器109Nを含む。プライマリノード101の中のこれらモジュールは、同期マッピング115、飛行中追跡データ(in-flight tracking data)117、及び飛行中追跡データ119として図1に表されるデータにアクセスする。同期マッピング115は、論理記憶オブジェクト間における同期構成を示す(本明細書において、同期関係(synchronization relationships)又は同期関係(sync relationships)ともいわれる)。例えば、プライマリ論理記憶オブジェクトが、1つのセカンダリ論理記憶オブジェクトとの完全同期関係と、別のセカンダリ論理記憶オブジェクトとの半同期関係を有することができる。飛行中追跡データは、要求の進捗又は状態を、対応する変更伝搬器の観点から追跡する。換言すると、各変更伝搬器インスタンスは、飛行中追跡データを、同期関係を有する対応した論理記憶オブジェクトについて維持する。

10

**【0016】**

セカンダリノード125は、プライマリノード101と同じモジュール/エンジンのすべてを含むことができる。図1において、モジュールのうちいくつかは、繰り返しを低減させるよう表されていない。セカンダリノード125は、ファイルシステム要求生成器129、記憶クラスタ同期エンジン139、変更伝搬エンジン131、及びファイルシステム137を含むものとして表されている。変更伝搬エンジン131は、傍受器133及びセカンダリライタ135を含む。セカンダリノード125のセカンダリライタ135は、同期追跡データ(sync tracking data)134として図1に表されるデータにアクセスする。同期追跡データ134は、要求の進捗又は状態を、セカンダリライタ135の観点から示す。同期追跡データ134は、必ずしも変更伝搬エンジン131内に含まれない。同期追跡データ134は、本説明のためにセカンダリライタ135の近くに単に表されている。セカンダリノード125は、論理記憶オブジェクトのグループ127をホストするホスト記憶要素に通信可能に結合される。グループ127は、セカンダリ記憶オブジェクトを含む。

20

30

**【0017】**

異なるノードにわたり機能性において様々なバリエーションがあり得るが、同じ名称を有するモジュールの機能性は一般に本例示において同じである。ファイルシステム要求生成器103、129は、ファイルシステム生成器103、129に渡される記憶プロトコル入力/出力(I/O)要求に基づいてファイルシステム要求を生成する。ファイルシステム生成器103、129は、記憶プロトコルI/O要求を、ネットワークスタック、スモールコンピュータシステムインターフェース(SCSI)スタック、インターネットSCSI(iSCSI)モジュール等から受信することができる。記憶プロトコルI/O要求の例には、ストレージエリアネットワーク(SAN)要求及びネットワークアタッチトストレージ(NAS)要求が含まれる。ファイルシステム生成器103、129は、そのノード上に実装されるファイルシステムに基づいてファイルシステム要求を生成する。傍受器107、133は、ファイルシステム要求生成器103、129からの要求を傍受する。傍受は、別様に実装されることができる。呼び出し元に提示されるインターフェースを変更することなく下層の機能性が変わるように、アプリケーションプログラミングインターフェースが修正されることができる。別の例として、監視処理が、実行キューを監視し、指定されたアドレスが実行キュー内で発生するときに呼び出しをリダイレクトすることができる。ファイルシステム111、137は、ファイルシステム要求に従って下層のホスト記憶要素にアクセスする。記憶クラスタ同期エンジン123、139は、ネットワーク110を介して実装されるプロトコルに従って通信を処理する。例として、エンジン123、139により実装されるプロトコルは、ファイバチャネル(FC)、ファイバチ

40

50

チャンネルオーバーイーサネット（登録商標）（FCoE）、インターネットファイバチャンネルプロトコル（iFCP）、及びトンネリングプロトコルのうち任意の1つ以上であり得る。特定のプロトコルにかかわらず、エンジン123、139は、マシン間の距離及びホップにかかわらずマシン間の直接接続として認識できるアクティブな接続をサポートするプロトコルを実装する。

#### 【0018】

図1は、文字A～Nによって識別される一連の段階（stages）を備えた例示的な動作を表す。文字によって示唆される動作の順序付けは本例示に限られ、請求項の範囲を限定するのに使用されるべきではない。段階Aにおいて、プライマリノード101が、変更要求113を受信する。変更要求113は、クライアントから発せられ、クライアントは、管理ノード（例えば、クラスタマネージャ）、ユーザノード（例えば、顧客のサーバ）等にあり得る。ファイルシステム要求生成器103が、変更要求113を処理し、変更要求に基づいてファイルシステム要求を生成し、ファイルシステム要求113をファイルシステム111に渡すコードを起動する。ファイルシステム要求113を生成することの一部として、ファイルシステム要求生成器103は、変更要求の中に対象として示される論理記憶オブジェクト識別子を、論理記憶オブジェクトのファイルシステム場所情報（例えば、inode識別子、オフセット等）に変換する。しかし、ファイルシステム要求生成器103は、ファイルシステム要求と共に進むように、上記論理記憶オブジェクト識別子をさらに示す。論理記憶オブジェクト識別子は、異なる仕方でもファイルシステム要求と共に進むことができる。例えば、ファイルシステム生成器は、論理オブジェクト識別子をファイルシステム要求のメタデータに書き込むことができる。別の例として、ファイルシステム生成器は、データ構造を作成し、それをファイルシステムに関連付ける。ファイルシステム111がファイルシステム要求を受信することに代わって、傍受器107が段階Bにおいてファイルシステム要求を受信する。様々な形式における要求（例えば、記憶プロトコルI/O要求、ファイルシステム要求等）はラベル113ではもはや識別されず、なぜならば、要求されている変更が要求の形式にかかわらず同じであるからである。

#### 【0019】

段階Cにおいて、傍受器107は、同期マッピング115にアクセスして、ファイルシステム要求に関連する何らかの同期関係を決定する。ファイルシステム要求は、（ファイルシステム場所情報の観点で）ファイルシステム要求の対象であるグループ102の中の論理記憶オブジェクトを示す。傍受器107は、同期マッピング115にアクセスして、ファイルシステム要求対象について定義されたいくらかの同期関係を決定する。対象は、単一の同期関係を有してもよく、複数の同期関係を有してもよく、あるいは同期関係を有さなくてもよい。対象が同期関係を有さない場合、ファイルシステム要求はファイルシステム111に渡されることになる。本例示について、同期マッピング115は、対象がグループ127の中の論理記憶オブジェクトとの完全同期関係を有することを示す。ファイルシステム要求の対象が同期関係を有するため、ファイルシステム要求の対象はプライマリ論理記憶オブジェクトとみなされることができる。前述されたとおり、論理記憶オブジェクトは、少なくとも互いに関連付けられたクラスタにわたり排他的である不変識別子によって識別される。同期マッピングは、1つ以上のデータ構造内に示されることがあり、同期マッピングは、論理オブジェクト（例えば、ファイル、LUN等）及び下層のファイルシステムに依存して、論理オブジェクト又はファイルシステム要求対象の複数レベル又は複数レイヤにわたり同期関係をマッピングする。例えば、論理オブジェクトはファイルであり得る。論理オブジェクト識別子は、最初、ファイル識別子又はファイルハンドルであることになる。ファイルシステムは、ファイルハンドルに影響されるデータブロックに向けて書き込み要求を解決する。ファイルシステムは、例えば、任意数のinodeレベルをとおして解決してもよい。同期関係が存在するとき、同期マッピングは、プライマリノードにおけるより高いレベルの識別子（すなわち、論理オブジェクト識別子）をセカンダリノードにおけるより高いレベルの識別子にマッピングするだけでなく、さらに、同期マッピングは、より低いレベルの識別子（すなわち、ファイルシステム場所情報）をマッ

10

20

30

40

50

ピングする。この例の場合、より低いレベルの識別子は、`inode` 識別子であることになる。対象にされているファイルの一部のためのプライマリノード `inode` 識別子は、対象にされているファイルの一部のためのセカンダリノード上の `inode` 識別子にマッピングすることになる。

#### 【0020】

段階Dにおいて、傍受器107は、ファイルシステム要求と対象に関する同期関係の指標とを変更伝搬器109Aに渡す。プライマリノード101が、変更要求113の中に示されるのと同じプライマリ論理記憶オブジェクトを対象にする変更要求をまだ受信していなかった場合、傍受器107は、変更伝搬器109Aをインスタンス化するコードを起動してもよい。必要ではないが、変更伝搬器は、この例示においてプライマリ論理記憶オブジェクトごとにインスタンス化される。傍受器107は、プライマリ論理記憶オブジェクトの同期関係を変更伝搬器に対して様々な仕方で示すことができる。例えば、傍受器107は、パラメータ値としてのプライマリ論理記憶オブジェクト識別子とパラメータ値としてのセカンダリ論理記憶オブジェクト識別子とを用いて、変更伝搬器をインスタンス化する関数を呼び出すことができる。別の例として、傍受器107は、すでにインスタンス化された変更伝搬器109Aに対して、ローカルメモリに記憶されたファイルシステム要求への参照と共に、プロセス間通信を送信することができる。プライマリ論理記憶オブジェクトごとの変更伝搬器のインスタンス化を例示するために、変更伝搬器109Nは、飛行中追跡データ119に対する破線と共に表されている。破線が使用されて、変更伝搬器109Nが異なるファイルシステム要求のための飛行中追跡データ119にアクセスしている可能性があることを示す。

10

20

#### 【0021】

段階Eにおいて、変更伝搬器109Aは、同期関係のセカンダリ論理記憶オブジェクトを対象にするファイルシステム要求を作成し、飛行中追跡データ117を更新する。変更伝搬器109Aがちょうどインスタンス化されたところである場合、追跡データのための構造がまだ存在しない可能性があり、あるいは空の構造が存在する可能性がある。変更伝搬器109Aは、飛行中追跡データ117を更新して、プライマリ論理記憶オブジェクトを対象にするファイルシステム要求が飛行中である（すなわち、送信されることになるか又は送信されている）ことを示す。変更伝搬器109Aは、飛行中追跡データ117を更新して、セカンダリ論理記憶オブジェクトを対象にするファイルシステム要求が飛行中であることをさらに示す。変更伝搬器109Aは、それから（又は同時に）、プライマリ論理記憶オブジェクトとの完全同期関係を有するセカンダリ論理記憶オブジェクトの識別子を備えた要求を作成する。変更伝搬器109Aは、異なる要求元を有するファイルシステム要求を同様に作成する。変更伝搬器109Aは、変更伝搬器109Aを要求元として示す。変更伝搬器109Aは、任意の関連付けられたクラスタ内で変更伝搬器109Aを排他的に識別する様々なデータを用いて識別されることができ、例えば、変更伝搬器109Aのプロセス/スレッド識別子とプライマリノード101のネットワークアドレスとの組み合わせなどである。変更伝搬器109Aは、プライマリ論理記憶オブジェクト識別子を要求元の指標にさらに組み込むことができる。変更伝搬器109Aから送信されるプライマリ論理記憶オブジェクトを対象にするファイルシステム要求は、プライマリ変更要求といわれる。変更伝搬器109Aから送信されるセカンダリ論理記憶オブジェクトを対象にするファイルシステム要求は、セカンダリ変更要求といわれる。

30

40

#### 【0022】

段階Fにおいて、変更伝搬器109Aは、サービス提供のためのファイルシステム要求を送信する。プライマリ論理記憶オブジェクトがセカンダリ論理記憶オブジェクトとの完全同期関係を有するため、プライマリノード101は、変更がプライマリ及びセカンダリ双方の論理記憶オブジェクトにおいて行われるまで、変更要求113に回答しないことになる。したがって、変更伝搬器109Aは、プライマリ及びセカンダリ変更要求を任意の順序で送信することができる。変更伝搬器109Aは、プライマリ変更要求をファイルシステム111に送信する。変更伝搬器109Aは、セカンダリ変更要求を記憶クラスタ同

50

期エンジン 1 2 3 に送信する。変更要求が変更伝搬器 1 0 9 A から渡された後、動作のタイミングは、ネットワーク条件、ノード能力における差等に依存して変動する可能性がある。

#### 【 0 0 2 3 】

段階 G において、ファイルシステム 1 1 1 は、ホスト記憶要素にアクセスする。段階 H において、記憶クラスタ同期エンジン 1 2 3 は、ネットワーク 1 1 0 を横断する記憶クラスタ同期エンジン 1 2 3 と記憶クラスタ同期エンジン 1 3 9 との間の接続のプロトコルに従って、セカンダリ変更要求を処理する。記憶クラスタ同期エンジン 1 2 3 は、接続プロトコルに従って新しい要求を構築し、この新しい要求をセカンダリ変更要求からの関連情報（例えば、セカンダリ論理記憶オブジェクト識別子、書き込まれるべきデータ等）で埋めることができる。記憶クラスタ同期エンジン 1 2 3 は、セカンダリ変更要求を、接続プロトコルに準拠したヘッダでカプセル化してもよい。この例示について、プライマリノードにおける同期マッピングは、プライマリノードとセカンダリノードとの間で論理オブジェクト識別子（例えば、ファイルハンドル）をマッピングし、さらに、ファイルシステム場所情報（例えば、i n o d e 識別子）をマッピングする。セカンダリ変更要求は、変更要求によって影響されるデータブロックのセカンダリノードファイルシステム場所情報を用いて構築される。いくつかの場合、ファイルシステム場所情報同期マッピングは、論理オブジェクト識別子同期マッピングとは別個であることになる。また、ファイルシステム場所情報同期マッピングは、セカンダリノードにおいて維持されてもよい。こうした場合、セカンダリ変更要求は、対象にされた論理オブジェクトの指標とプライマリノードのファイルシステム場所情報とを用いて構築される。受信されたとき、セカンダリノードは、同期マッピングにアクセスし、セカンダリノードファイルシステム場所情報に対するプライマリノードファイルシステム場所情報を解決することになる。

10

20

#### 【 0 0 2 4 】

段階 I において、記憶クラスタ同期エンジン 1 3 9 が、接続プロトコルに従って受信した要求を処理し、セカンダリ変更要求をセカンダリライター 1 3 5 に渡す。記憶クラスタ同期エンジン 1 3 9 は、受信した要求からセカンダリ変更要求を再構築し、あるいは受信した要求からセカンダリ変更要求を抽出することができる。セカンダリ変更要求がまだ受信されていなかった場合、記憶クラスタ同期エンジン 1 3 9 は、セカンダリライター 1 3 5 をインスタンス化するコードを起動してもよい。記憶クラスタ同期エンジン 1 3 9 は、記憶クラスタ同期エンジン 1 3 9 により受信されるすべてのセカンダリ変更要求を扱うセカンダリライターをインスタンス化し、あるいはプライマリ論理記憶オブジェクト及びセカンダリ論理記憶オブジェクトのペアごとにこれらをインスタンス化することができる。

30

#### 【 0 0 2 5 】

図 1 は、ファイルシステム要求生成器 1 2 9 から及び傍受器 1 3 3 からの破線を表している。ファイルシステム要求生成器 1 2 9 からの破線は、ファイルシステム要求生成器 1 2 9 が他の変更要求を受信し、処理し、傍受器 1 3 3 に渡している可能性を示す。傍受器 1 3 3 から省略記号への破線は、傍受器 1 3 3 が変更要求を傍受し、セカンダリノード 1 2 5 の変更伝搬器に渡している可能性を例示しており、該変更伝搬器は図示されていない。こうした可能性は、セカンダリノード 1 2 5 がセカンダリ変更要求を扱うことに限定されないことを示すために例示される。

40

#### 【 0 0 2 6 】

段階 J において、セカンダリライター 1 3 5 が、同期追跡データ 1 3 4 を更新する。セカンダリライター 1 3 5 は、対象にされたセカンダリ論理記憶オブジェクトと要求元（すなわち、変更伝搬器 1 0 9 A）とセカンダリ変更要求の状態とを少なくとも含むセカンダリ変更要求の指標を記録する。この時点で、セカンダリライター 1 3 5 は状態を飛行中として記録し、なぜならば、セカンダリ変更要求は送信されているか又は送信されることになるからである。段階 K において、セカンダリライター 1 3 5 は、セカンダリ変更要求をファイルシステム 1 3 7 に送信する。

#### 【 0 0 2 7 】

50

段階 L において、ファイルシステム 137 が、セカンダリ変更要求に従ってホスト記憶要素にアクセスする。

【0028】

図 2 は、図 1 の同期マッピングにおいて定義される完全同期関係に従って処理されるプライマリ及びセカンダリ変更要求に対する応答を表す。図 2 は、段階ラベル A ~ L を用いて例示的な動作を表す。段階 A ~ J は、セカンダリノード 125 の前にプライマリ論理記憶オブジェクトのホスト記憶要素からの応答が応答する場合として表されている。しかしながら、その順序付けは必要ではない。いくつかの場合、プライマリ論理記憶オブジェクトのホスト記憶要素がプライマリノード 101 に応答することができる前に、セカンダリノード 125 が変更伝搬器 109A に応答することができる可能性がある。応答のタイミングにかかわらず、プライマリ及びセカンダリ双方の論理記憶オブジェクトにおける変更が変更伝搬器 109A によって確認されるまで、要求元に対する応答は提供されない。図 1 からのいくつかの要素が、図 2 を簡素化するために除去されている。

10

【0029】

段階 A ~ C は、プライマリ論理記憶オブジェクトのホスト記憶要素から変更伝搬器 109A に進む応答と、飛行中追跡データ 117 の対応する更新とを例示する。段階 A において、プライマリ論理記憶オブジェクトをホストするホスト記憶要素が、ファイルシステム 111 に応答を供給する。ファイルシステム 111 は、段階 B においてこの応答を変更伝搬器 109A に転送する。段階 C において、変更伝搬器 109A は、飛行中追跡データ 117 を更新して、プライマリ変更要求がプライマリ論理記憶オブジェクトにおいて実行されたことを示す。

20

【0030】

段階 D ~ J は、セカンダリ論理記憶オブジェクトのホスト記憶要素から変更伝搬器 109A に進む応答と、飛行中追跡データ 117 の対応する更新とを例示する。段階 D において、セカンダリ論理記憶オブジェクトをホストするホスト記憶要素が、応答をファイルシステム 137 に供給する。ファイルシステム 137 は、段階 E においてこの応答をセカンダリライタ 135 に転送する。段階 F において、セカンダリライタ 135 は、同期追跡データを更新して、セカンダリ論理記憶オブジェクトに対する更新を反映する。例えば、セカンダリライタ 135 は、セカンダリ論理記憶オブジェクト識別子と転送される応答の要求元との組み合わせを使用して、同期追跡データ 134 をホストする構造の中のエントリをロックアップする。セカンダリライタ 135 は、上記エントリの中に値又はフラグを設定して、セカンダリ論理記憶オブジェクトに対して変更が完了したことを示す。セカンダリライタ 135 は、それから、応答を記憶クラスタ同期エンジン 139 に転送する。記憶クラスタ同期エンジン 139 は、セカンダリ変更要求に対する応答（「セカンダリ応答」）がプライマリノード 101 に対して送信されるべきかを決定する。記憶クラスタ同期エンジン 139 は、接続プロトコルに従ってセカンダリ応答を処理し、段階 H においてネットワーク 110 を介して接続を通じてセカンダリ応答を送信する。段階 I において、記憶クラスタ同期エンジン 123 が、接続プロトコルに従ってセカンダリ応答を処理し、セカンダリ応答を変更伝搬器 109A に転送する。セカンダリ応答の処理の一部として、記憶クラスタ同期エンジン 123 は、変更伝搬器 109A のプロセス/スレッド識別子を組み込んである要求元識別子に基づいて、上記セカンダリ応答が変更伝搬器 A に送信されるべきかを決定することができる。段階 J において、変更伝搬器 109A は、飛行中追跡データ 117 を更新して、セカンダリ変更要求がセカンダリ論理記憶オブジェクトにおいて実行されたことを示す。

30

40

【0031】

最初の変更要求 113 に対応するすべての未処理の変更要求が完了したと決定した後、変更伝搬器 109A は、応答をファイルシステム要求生成器 103 に供給する。変更伝搬器 109A が飛行中追跡データ 117 を更新するたび、変更伝搬器 109A は、例えば、エントリを読み出して、エントリ内に示されるすべての要求が完了したか又は依然として飛行中であるかを決定することができる。この例示について、ファイルシステム要求生成

50

器 1 0 3 は、変更要求 1 1 3 に対応する要求元を示すデータを維持する。要求がファイルシステム要求生成器 1 0 3 によって最初に受信されるとき、この要求は、要求元に対応する要求識別子でタグ付けされることができる。この要求識別子は、要求及び対応する応答と共に進むことができる。要求識別子は、要求元のアイデンティティと要求とを示して、該要求を同じ要求元からの他の要求から区別する。変更伝搬エンジン 1 0 5 は、変更要求 1 1 3 の要求元を示し及び変更要求 1 1 3 それ自体を示すデータを、さらに（又は代わって）維持するようにプログラムされることができる。段階 L において、ファイルシステム要求生成器 1 0 3 は、変更応答 2 1 3 を形成し、変更応答 2 1 3 を対応する要求元に供給する。

#### 【 0 0 3 2 】

論理記憶オブジェクト間でとり得る同期関係の組み合わせのさらなる一例示として、図 3 ~ 4 は、完全同期関係においてセカンダリ論理記憶オブジェクトとして構成され及び半同期関係においてプライマリ論理記憶オブジェクトとして構成された論理記憶オブジェクトについての例示的な動作を表す。論理記憶オブジェクトの異なる観点を提供するために、図 3 ~ 4 は、ホスト記憶要素のクラスタの文脈において論理記憶オブジェクトを表す。論理記憶オブジェクトは、ホスト記憶要素（例えば、記憶アレイ）上に破線を用いて表される。論理記憶オブジェクトがこのように表されて、論理記憶オブジェクトが複数のホスト記憶要素に及ぶことと、単一のホスト記憶要素内にホストされることとの可能性を例示する。ホスト記憶要素が、ホスト記憶要素の集合（例えば、ディスクアレイ）である場合、論理記憶オブジェクトは、ディスクアレイ内の複数ディスクに及ぶ可能性がある。図 3 は、ノード 3 0 1 に関連付けられた記憶クラスタ 3 0 3 を表す。図 3 ~ 4 は、ノード 3 1 1 に関連付けられた記憶クラスタ 3 2 5 と、ノード 3 2 9 に関連付けられた記憶クラスタ 3 3 1 とを表す。図 3 ~ 4 は、ネットワーク 3 0 9 を介して通信するノードを表す。ノード 3 0 1 は、図 1 のノード 1 0 1 と同様に動作し、ゆえにその動作は、例示詳細について、図 1 内と同じレベルでは表されない。同様に、ノード 3 2 9 は、図 1 及び 2 のセカンダリノード 1 2 5 と同様に動作し、ゆえにその例示的な動作もまた、この例示に対して全体では繰り返されない。

#### 【 0 0 3 3 】

図 3 ~ 4 は、図 1 に表されるモジュールのうちいくつかを有するノード 3 1 1 を表す。再びになるが、繰り返しを回避するため、モジュールのすべてが繰り返されてはいない。図 3 ~ 4 において、N O D E \_ 2 として識別されるノード 3 1 1 が、セカンダリライタ 3 1 5、ファイルシステム 3 2 1、及び記憶クラスタ同期エンジン 3 1 3 を含む。図 3 ~ 4 は、例示的な同期関係を有する同期マッピング 3 1 7 としてノード 3 1 1 内に同期マッピングデータをさらに表している。さらに、ノード 3 1 1 は、追跡データ 3 1 9 を有する。しかし、追跡データ 3 1 9 は、ノード 3 1 1 からホスト記憶要素に送信される要求の状態と、ノード 3 1 1 から別のノードに送信される要求の状態とを示す。追跡データ 3 1 9 は、図 1 ~ 2 の飛行中追跡データと同様である。図 1 ~ 2 における変更伝搬器及びセカンダリライタの表現とは違って、図 3 ~ 4 は、変更伝搬器を、別の変更伝搬器からの変更要求に応答し、別の場所におけるセカンダリ論理記憶オブジェクトに変更を伝搬する機能を有するものとして表している。図 3 ~ 4 は、オブジェクト場所データ 3 2 7 をさらに表す。異なる例示的なエントリが図 3 ~ 4 に表されているが、こうしたエントリがオブジェクト場所データ内にあり、ノード識別子に対する論理記憶オブジェクト識別子を解決する。図 1 ~ 2 内と同様に、図 3 ~ 4 内の段階は、段階識別子と共に例示的な動作を表す。段階識別子は動作におけるシーケンスを示すが、その表される順序は請求項の範囲を限定するのに使用されるべきではなく、なぜならば、上記順序は例示目的のものだからである。

#### 【 0 0 3 4 】

段階 A ~ C は、図 1 における段階 C、H、及び I と同様である。段階 A において、ノード 3 0 1 は、変更要求を受信した後、同期マッピング 3 0 5 にアクセスする。該変更要求は、表されていない。同期マッピング 3 0 5 を用いて、ノード 3 0 1 は、O B J 3 3 として識別される論理記憶オブジェクトが O B J 4 4 として識別されるオブジェクトとの完全

10

20

30

40

50

同期関係を有すると決定する。論理記憶オブジェクトOBJ33は、上記関係においてプライマリ論理記憶オブジェクトであり、記憶クラスタ303内にホストされている。記憶クラスタ303は、ノード301に関連付けられている。ノード301は、段階Bにおいて、OBJ33をホストする記憶クラスタ303のメンバに変更要求を送信する。段階Cにおいて、ノード301は、オブジェクト場所データ307にアクセスし、OBJ44がNODE\_\_2とNODE\_\_2のアドレスとに関連付けられていることを決定する。NODE\_\_2はノード311である。ノード301は、それから、セカンダリ変更要求を送信する。セカンダリ変更要求は、(ファイルシステムの観点での)対象としてのOBJ44と、要求元としてのノード301内の変更伝搬器とを、ネットワーク309を介したノード間の接続を通じてノード311に対して示す。

10

**【0035】**

段階D~Gにおいて、ノード311は、ノード301からのセカンダリ要求を処理する。段階Dにおいて、記憶クラスタ同期エンジン313は、接続のプロトコルに従ってノード301からのセカンダリ要求を処理する。記憶クラスタ同期エンジン313は、それから、セカンダリ変更要求をセカンダリライタ315に渡す。セカンダリライタ315は、段階Eにおいて同期マッピング317にアクセスする。セカンダリライタ315は、論理記憶オブジェクトOBJ44が論理記憶オブジェクトOBJ52との半同期関係を有すると決定する。段階Fにおいて、セカンダリライタ315は、追跡データ319を更新する。セカンダリライタ315は、決定された半同期関係に基づいて作成されることになるセカンダリ変更要求について、及びノード301から受信されるセカンダリ変更要求について、追跡データ319を更新する。セカンダリライタ315は、追跡データ319内に、応答をどこにルーティングするかを指標を維持する。この例において、セカンダリライタ315は、追跡データ319を更新して、オブジェクトOBJ44及びOBJ52が半同期関係にあることを示す。論理記憶オブジェクト識別子の各々が、状態指標に関連付けられる。この例示について、“0”の値は、飛行中又は待ちを示し、“1”の値は、対象論理記憶オブジェクトに対して変更要求が実行されたことを示す。この時点で、双方の状態指標は“0”に設定される。セカンダリライタ315は、追跡データ319をさらに更新して、要求元を“NODE\_\_1\_\_OBJ33”として示す。この値は単に、同期関係のノード及びプライマリ論理記憶オブジェクトの一例示的な指標である。段階Gにおいて、セカンダリライタ315は、さらなるセカンダリ変更要求を作成し、変更要求をその対応するハンドラに渡す。セカンダリライタ315は、OBJ44を対象にする変更要求をファイルシステム321に転送する。変更伝搬器315は、OBJ52としての対象と、セカンダリライタ315及びノード311を示す要求元識別子とを用いて、さらなる変更要求を作成する。例えば、さらなる要求は、ノード311と、セカンダリライタ315に結び付けられたポート又はソケットとを示してもよい。セカンダリライタ315は、さらなる変更要求を記憶クラスタ同期エンジン313に渡す。

20

30

**【0036】**

段階Hにおいて、ファイルシステム321は、セカンダリ要求に従って、OBJ44をホストするクラスタ325内のホスト記憶要素にアクセスする。

**【0037】**

段階Iにおいて、記憶クラスタ同期エンジン313は、セカンダリライタ315からのさらなる変更要求をどこに送信するかを決定する。記憶クラスタ同期エンジン313は、オブジェクト場所データ327にアクセスし、OBJ52がNODE\_\_3に関連付けられていることを示すエントリを発見する。NODE\_\_3はノード329である。記憶クラスタ同期エンジン313は、オブジェクト場所データ327からノード329のアドレスを決定し、接続プロトコルに従ってさらなる変更要求を処理し、段階Jにおいてさらなる変更要求をネットワーク309を介してノード329に送信する。段階Kにおいて、ノード329は、OBJ52に対してさらなる変更要求を実行する。

40

**【0038】**

図4は、ノード311による、異なる同期関係のための応答の扱いを表す。段階A~D

50



は、ノード 3 1 1 が O B J 4 4 を対象にした変更要求に対する応答を処理する例示的な動作を表す。段階 A において、O B J 4 4 をホストするホスト記憶要素が、応答をファイルシステム 3 2 1 に送信する。段階 B において、ファイルシステム 3 2 1 が応答をセカンダリライタ 3 1 5 に転送し、なぜならば、セカンダリライタ 3 1 5 が変更要求の発信元 (originator) であったことを上記応答が示すからである。段階 C において、セカンダリライタ 3 1 5 は、追跡データ 3 1 9 にアクセスする。セカンダリライタ 3 1 5 は、追跡データ 3 1 9 を更新して、O B J 4 4 に対する変更が実行されたことを示す。セカンダリライタ 3 1 5 は、O B J 4 4 のためのエントリが半同期関係を示し、要求元 N O D E \_ 1 \_ O B J 3 3 を示すと決定する。上記が半同期関係であるため、変更伝搬器 3 1 5 は、要求元 N O D E \_ 1 \_ O B J 4 4 に対する応答の提供を進めることができる。セカンダリライタ 3 1 5 は、段階 D において応答を要求元識別子の指標と共に記憶クラスタ同期エンジン 3 1 3 に送信する。

10

20

30

40

50

**【 0 0 3 9 】**

O B J 3 3 と O B J 4 4 との間の同期関係は完全同期関係であるため、O B J 4 4 に対する変更は、O B J 3 3 に関連付けられたノードに逆向きに (back) 即座に通信されることができる。段階 E において、記憶クラスタ同期エンジン 3 1 3 は、オブジェクト場所データ 3 2 7 にアクセスして、O B J 3 3 に関連付けられたノードを決定する。記憶クラスタ同期エンジン 3 1 3 は、セカンダリライタ 3 1 5 によって提供される要求元識別子からオブジェクト識別子を抽出するようにプログラムされることができる。しかしながら、オブジェクト識別子は異なる仕方で通信されてもよい。例えば、セカンダリライタ 3 1 5 又は記憶クラスタ同期エンジン 3 1 3 が同期マッピングにアクセスして、O B J 3 3 のためのプライマリ論理記憶オブジェクトを決定することができる。オブジェクト識別子が如何にして決定されるかにかかわらず、記憶クラスタ同期エンジン 3 1 3 は、O B J 3 3 が N O D E \_ 1 に関連付けられていると決定する。N O D E \_ 1 はノード 3 0 1 である。

**【 0 0 4 0 】**

段階 F において、応答がノード 3 0 1 に送信される。記憶クラスタ同期エンジン 3 1 3 は、段階 E において逆向きに宛先を N O D E \_ 1 として決定した後 (又は決定する間)、接続プロトコルに従って応答を処理する。記憶クラスタ同期エンジン 3 1 3 は、それから、ネットワーク 3 0 9 を横断する接続をとおして応答を伝達する。ノード 3 0 1 は、それから、応答を作成し、ネットワーク 4 0 1 を通じて最初の要求 4 0 3 (「クライアント」) に送信する。このことは、要求される変更が O B J 3 3 においてすでに実行されていることを仮定し、なぜならば、O B J 3 3 が O B J 4 4 との完全同期関係を有するからである。

**【 0 0 4 1 】**

段階 H ~ J において、O B J 5 2 に対する変更を確認する応答が、逆向きにセカンダリライタ 3 1 5 に進む。段階 H において、O B J 5 2 をホストする記憶クラスタ 3 3 1 のメンバが、ノード 3 2 9 に、O B J 5 2 に対して変更が実行されたとの応答を提供する。したがって、ノード 3 2 9 は、段階 I において、応答を記憶クラスタ同期エンジン 3 1 3 に送信する。応答がセカンダリライタ 3 1 5 を示すので、記憶クラスタ同期エンジン 3 1 3 は、接続プロトコルに従って応答を処理した後、段階 J においてセカンダリライタ 3 1 5 に応答を渡す。

**【 0 0 4 2 】**

段階 K において、変更伝搬器 3 1 5 は、追跡データ 3 1 9 を更新して、O B J 5 2 に対する更新が完了したことを示す。同期が完了したとのこの指標は、一貫性についての他の態様に対して使用されることができ、例えば、シーケンシング、フェイルオーバー、及びロードバランシングなどである。

**【 0 0 4 3 】**

図 1 ~ 4 はアーキテクチャの一部を表して例示的な動作を示しているが、図 5 は、1 つ又は複数のクラスタのノードにわたり論理記憶オブジェクト粒度において完全同期と半同期との双方を提供する一例示的なアーキテクチャを表す。図 5 は、ファイルシステム要求

生成器 501、変更伝搬エンジン 503、ファイルシステム 505、及び記憶クラスタ同期エンジン 507 を表している。ファイルシステム要求生成器 501 は、図 1 ~ 4 のファイルシステム要求生成器 103 と同様である。ファイルシステム要求生成器 501 は、ネットワークインターフェース又はシリアルインターフェース（例えば、ネットワークモジュール/スタック又は SCSI モジュール）を通じて受信される通信を処理するモジュールから受信される記憶プロトコル固有 I/O 要求を処理する。ファイルシステム 505 は、図 1 のファイルシステム 111 と同様とすることができ、ファイルシステム又はファイルシステムレイヤを実装する。ファイルシステム又はファイルシステムレイヤの例には、Write Anywhere File Layout 及び UNIX（登録商標）ファイルシステムが含まれる。ファイルシステム 505 は、ファイルシステム要求に従って下層のホスト記憶要素に要求を供給する。記憶クラスタ同期エンジン 507 は、同期関係のセカンダリ論理記憶オブジェクトをホストするクラスタノードにおけるカウンターパートの記憶クラスタ同期エンジンに、変更要求を供給する。

10

20

30

40

50

#### 【0044】

変更伝搬エンジン 503 は、傍受器 509、シーケンサ (sequencer) 511、伝搬器 513、及びセカンダリライタ 515 を含む。ファイルシステム要求生成器 501 は、ファイルシステム要求を、対応する記憶 I/O 要求（例えば、SAN 又は NAS 要求）内に示された論理記憶オブジェクト対象の指標と共に、変更伝搬エンジン 503 に渡す。ファイルシステム要求生成器 501 の観点から、ファイルシステム要求生成器 501 は、ファイルシステム要求をファイルシステム 505 に渡している。このことは、ファイルシステム要求生成器 501 の修正を回避し又は最小化することに役立つ可能性がある。しかし、ファイルシステム要求生成器 501 によって起動される関数又はプロシージャコールが傍受器 509 を実際には起動し、ゆえに、傍受器 509 がファイルシステム要求を「傍受する」ことが可能になる。応答が伝搬器 513 から受信されるとき、傍受器 509 は、応答を逆向きにファイルシステム要求生成器 501 に渡し、ファイルシステム要求生成器 501 は、それから、対応する記憶プロトコル I/O 応答を作成する。

#### 【0045】

ファイルシステム要求は、最初、傍受器 509 に渡る。傍受器 509 は、最初、ファイルシステム要求が如何にして変更伝搬エンジンを通して流れるかを決定する。ファイルシステム要求が変更要求（例えば、書込み、ゼロ設定 (zero) 等）である場合、傍受器 509 は、同期関係の中の論理記憶オブジェクトを示す同期関係データにアクセスする（同期関係を有する論理記憶オブジェクトは、以降、エンドポイントといわれる）。同期関係データが、変更要求の対象（すなわち、プライマリ論理記憶オブジェクトであり、以降、「プライマリエンドポイント」といわれる）とセカンダリ論理記憶オブジェクト（すなわち、プライマリエンドポイントと同期する論理記憶オブジェクトであり、以降、「セカンダリエンドポイント」といわれる）との間の完全同期関係を示す場合、傍受器 509 は、変更要求、同期関係の指標、及びセカンダリエンドポイントの指標を伝搬器 513 に渡す。傍受器 509 は、この情報を、伝搬器 513 のインスタンス化を結果としてもたらず関数を呼び出すことによって渡すことができる。同期関係が半同期関係である場合、傍受器 509 は、この情報をシーケンサ 511 に渡す。変更伝搬エンジン 503 は、ファイルシステム要求と対応する半同期関係情報とを、シーケンサ 511 と伝搬器 513 との双方に対して、同時に又は時間において互いに近接して渡す傍受器を備えて設計されることができる。伝搬器 513 と同様に、傍受器 509 は上記情報をシーケンサ 511 に関数コールを用いて渡すことができ、上記関数コールがシーケンサ 511 をインスタンス化することができる。シーケンサ 511 及び伝搬器 513 は、プライマリ及びセカンダリエンドポイントの各ペアについてインスタンス化される。

#### 【0046】

シーケンサ 511 は、半同期関係におけるエンドポイントのためのファイルシステム要求で、又は、特定の記憶管理動作がトリガされるとき、動作する。上記管理動作は、例えば、スナップショット又は重複排除 (deduplication) などである。シーケンサ 511 は

、依存関係 (dependencies) を有する要求の順序を保存する。依存関係は、重なった書き込み間で、書き込みの間の読み出し要求から、指定された依存関係から等で生じることがある。シーケンサ 5 1 1 は、RPO などの構成の制限内で変更要求を追跡する。例えば、RPO が 10 秒として定義されるとき、シーケンサが 5 秒間隔で要求を追跡してもよい。個別の構成が何であれ、シーケンサ 5 1 1 は、構成された境界の範囲内で変更要求を蓄積し、上記境界は、時間、要求の数、又は双方の観点におけるものとするができる。本説明は、有界の蓄積された変更要求を変更セットという。シーケンサ 5 1 1 は、蓄積された変更要求間における依存関係を決定し、依存関係に基づいてシーケンスを示す。シーケンサ 5 1 1 は、読み出し要求と変更要求とについての可視性を有して、変更要求間における依存関係を決定する。変更セットの中の各要求のメタデータの中に、シーケンサ 5 1 1 は、変更セット内におけるシーケンシングと要求の総数とを示す。例えば、シーケンサ 5 1 1 は、変更セット内の 5 つの変更要求のうちの第 1 の変更要求のためのメタデータを “1 / 5” として書き込む。変更セットの境界が到達されるとき、シーケンサ 5 1 1 は、次の変更セットのための要求を蓄積することと、セカンダリエンドポイントをホストするノードに対して通信するために同期エンジン 5 0 7 に現在の変更セットを通信することとを開始する。本説明は、上記処理を、現在の変更セットのクローズ又は現在の変更セットログのクローズ、及び、次の変更セット又は変更セットログのオープンという。シーケンサ 5 1 1 は、示された順序で変更セットログを横断する別のスレッド又はプロセス (例えば、バックグラウンドプロセス) を立ち上げることができ、各変更要求を同期エンジン 5 0 7 にサブミットする (submits)。シーケンサ 5 1 1 (又は、シーケンサ 5 1 1 により起動されるスレッド/プロセス) は、変更セットの中の変更要求を個々に送信し、このことは、要求が順序不同で受信ノードに到着することを可能にする。変更セットについて、成功の応答が同期エンジン 5 0 7 から受信されるとき、シーケンサ 5 1 1 は、この変更セットを完了としてマーク付けする。それから、変更セットログは破棄され、あるいは上書きされることができる。失敗の応答が受信され、あるいはタイムアウトが発生するとき、シーケンサ 5 1 1 は、同期が失敗したとの通知を生成し、あるいは再試行することができる。

10

20

30

40

50

#### 【0047】

伝搬器 5 1 3 は、変更要求の状態を追跡するデータを維持し、セカンダリエンドポイントのための要求を同期エンジン 5 0 7 に渡し、応答を逆向きに傍受器 5 0 9 に渡す。伝搬器 5 1 3 が傍受器 5 0 9 から変更要求を受信するとき、伝搬器 5 1 3 は、要求元の指標を記録し、それから、変更要求を修正して伝搬器 5 1 3 を要求元として示す。このことは、ファイルシステム 5 0 5 が伝搬器に応答を返すのを容易にするが、必要ではない。要求元のアイデンティティを変更することに代わって、ファイルシステムからの応答を傍受するようにアーキテクチャが設計されることができる。要求元のアイデンティティを変更するアーキテクチャにおいて、伝搬器 5 1 3 は、傍受器に応答を渡す前に、ファイルシステム 5 0 5 からの応答の中の要求元のアイデンティティを復元する。変更要求を扱うことに戻ると、伝搬器 5 1 3 は、完了していない変更要求 (すなわち、飛行中変更要求) を示すデータを記録する。完全同期について、伝搬器 5 1 3 は、プライマリエンドポイント及びセカンダリエンドポイントのための変更要求に関するデータを記録する。伝搬器 5 1 3 は、このデータを使用して、いつ双方が完了するかと応答が要求元に提供されることができるかとを決定する。半同期について、伝搬器 5 1 3 は、プライマリエンドポイントのためのデータを記録し、なぜならば、シーケンサ 5 1 1 が半同期関係においてセカンダリエンドポイントのための要求を扱うからである。しかし、伝搬器 5 1 3 は、シーケンサ 5 1 1 に、プライマリエンドポイント上でいつ変更が完了されるかを通知する。こうした変更のすべてがプライマリエンドポイント上で成功裏に完了するまで、シーケンサ 5 1 1 はセカンダリエンドポイントのための変更セットを送出ししない。

#### 【0048】

セカンダリライタ 5 1 5 は、セカンダリエンドポイントを対象にする変更要求を扱う。セカンダリライタ 5 1 5 は、変更要求を複製動作の形式で同期エンジン 5 0 7 から受信す

る。プライマリエンドポイントノードにおいて、同期エンジン 5 0 7 は、伝搬器 5 1 3 又はシーケンサ/シーケンサによりスポンされた (spawned) スレッド 5 1 1 から供給される変更要求から、複製動作を生成する。複製動作は、プライマリエンドポイントノード (例えば、プライマリエンドポイントノード上の伝搬器インスタンス) を複製動作のソースとして示し、セカンダリエンドポイントを示す。複製動作は、同期関係のタイプをさらに示す。複製動作は、変更要求とは異なる要求元及び異なる対象を示し、ファイルシステム 5 0 5 のプロトコルから独立したプロトコルに従う、再形成された変更要求であってもよい。例えば、伝搬器からの変更要求が、要求元として伝搬器を、並びに、W A F L (w r i t e a n y w h e r e f i l e l a y o u t) に準拠する要求の中の対象として個別のファイル名及びファイル領域 (例えば、ブロック) を示してもよい。複製動作は、上記情報を変更要求から抽出し、これを個別のファイルシステム又はプロトコルから独立して示すことができる。セカンダリライタ 5 1 5 が上記情報を受信するとき、セカンダリライタ 5 1 5 は、セカンダリエンドポイントノードのファイルシステムにより実装されるプロトコルに従って適切な要求を生成する。セカンダリライタ 5 1 5 は、ファイルシステム 5 0 5 にサブミットされる要求の状態を追跡し、応答を逆向きに同期エンジン 5 0 7 に渡す。半同期関係について、セカンダリライタ 5 1 5 は、変更セットが完了するまで、変更セットの要求を蓄積する (「ステージする (stages) 」)。セカンダリライタ 5 1 5 は、要求のメタデータを読み出して、変更セットをいつ作成するかと変更セットがいつ完了するかとを決定する。変更セットが完了するとき、セカンダリライタ 5 1 5 は、プライマリエンドポイントノードに対して、変更セット内の各要求について個々の通知を送信することに代わって、変更セットが成功裏に完了したとの通知を生成することになる。セカンダリライタ 5 1 5 は、失敗の変更セットについて、通知をさらに生成することになる。いくつかの場合、セカンダリエンドポイントは、別の同期関係におけるプライマリエンドポイントである (「カスケード的 (cascading) 同期構成」)。セカンダリライタ 5 1 5 がインスタンス化されるとき、このセカンダリライタインスタンスは、同期関係データにアクセスして、セカンダリエンドポイントがカスケード的同期構成にあるかを決定する。そうである場合、セカンダリライタ 5 1 5 は、カスケード的關係のために伝搬器インスタンス及び / 又はシーケンサインスタンスを起動することになる。セカンダリライタ 5 1 5 は、それ自体が変更要求の要求元を有すると示すことになる。

10

20

30

#### 【 0 0 4 9 】

図 5 は、一例示的な論理オブジェクト粒度完全同期及び半同期アーキテクチャの一般的説明を提供するが、以降の図は、例示的な動作についてのさらなる例示を提供する。図 6 ~ 1 3 は、クラスタノード内のエンドポイントについての例示的な完全同期及び半同期動作のフローチャートを表す。これら図は、図 5 に表される例示的なアーキテクチャからのアクタを参照して説明されるが、指定されるアクタは動作の理解を助けるためのものである。前述されたとおり、プログラム構造又は設計が変動する可能性があり、アクタを指定する例は請求項の範囲を限定するのに使用されるべきではない。

#### 【 0 0 5 0 】

図 6 は、変更要求の受信を扱い、完全同期関係におけるプライマリエンドポイントを対象にする変更要求を扱う、例示的な動作のフローチャートを表す。傍受器が、ブロック 6 0 1、6 0 3、6 0 5、及び 6 0 7 の動作を実行することができ、伝搬器が、ブロック 6 0 8、6 0 9、6 1 1、及び 6 1 3 の動作を実行することができる。

40

#### 【 0 0 5 1 】

ブロック 6 0 1 において、傍受器が、記憶プロトコル I / O 要求から導出されるファイルシステム要求 (以降、「要求」) を受信する。例えば、情報が記憶プロトコル I / O 要求から抽出されて、ファイルシステム要求が生成されている。

#### 【 0 0 5 2 】

ブロック 6 0 3 において、傍受器は、要求の対象が同期関係にあるかどうかを決定する。要求の対象が同期関係にない場合、制御はブロック 6 1 7 に流れる。要求の対象が同期関係にある場合、制御はブロック 6 0 4 に流れる。

50

## 【 0 0 5 3 】

ブロック 6 0 4 において、傍受器は、同期関係情報を読み出す。傍受器は、情報を別のデータ構造にコピーすることと、そのデータ構造を変更要求に関連付けることとによって、上記情報を「読み出す」ことができる。傍受器は、さらに、情報を含むデータ構造の中のエントリに対する参照（例えば、ポインタ、索引等）を記録することによって、上記情報を「読み出す」ことができる。同期関係情報は、同期関係におけるエンドポイントをホストする（1つ以上の）クラスタのメンバ間で広められる（circulated）データ構造の中に維持されることができる。このデータ構造は、個々のノードにおいて構成されることができる。構成は、コミットされた後、（1つ以上の）クラスタ内のノードにわたる更新をトリガすることができる。

10

## 【 0 0 5 4 】

ブロック 6 0 5 において、傍受器は、要求が変更要求又は読み出し要求であるかを決定する。傍受器が、要求が変更要求であると決定する場合、制御はブロック 6 1 1 に流れる。そうでなければ、制御はブロック 6 0 7 に流れる。

## 【 0 0 5 5 】

ブロック 6 0 7 において、傍受器は、要求の同期関係が完全同期関係又は半同期関係であるかを決定する。関係が完全同期関係である場合、制御はブロック 6 1 7 に流れ、なぜならば、読み出しは、対応する完全同期動作をトリガしないからである。関係が半同期関係である場合、制御はブロック 6 0 9 に流れ、なぜならば、読み出しが、変更セットの中の変更要求間における依存関係を作り出す可能性があるからである。

20

## 【 0 0 5 6 】

ブロック 6 0 9 において、要求は、同期関係のプライマリエンドポイント及びセカンダリエンドポイントペアのためのシーケンサインスタンスに渡される。シーケンサは、上記ペアについてすでにインスタンス化されていてもよく、あるいは、要求を渡すことと同時にインスタンス化されてもよい。例えば、傍受器が、各々の一意のエンドポイントペアリングについてインスタンス化されたシーケンサを示すデータをチェックすることができる。傍受器がエントリを発見する場合、傍受器は、変更要求及び同期関係情報に対する参照を、上記エントリの中のスレッド識別子を用いてスレッドに渡す。エントリが存在しない場合、傍受器は関数を呼び出し、変更要求及び同期関係情報に対する参照が該関数コールのパラメータとして渡される。制御はブロック 6 0 9 から図 8 のブロック 8 0 1 に流れる。

30

## 【 0 0 5 7 】

傍受器がブロック 6 0 5 において、要求が変更要求であったと決定した場合、制御はブロック 6 1 1 に流れている。ブロック 6 1 1 において、要求及び同期関係情報が、同期関係情報の中に示されるプライマリエンドポイント及びセカンダリエンドポイントペアのための伝搬器インスタンスに渡される。シーケンサと同様に、伝搬器インスタンスは、要求及び同期関係情報を渡すことと同時にインスタンス化されてもよい。さらに、変更要求と同期関係情報とのうちいずれか又は双方が、参照で（referentially）又はリテラルで（literally）渡されることができる。

40

## 【 0 0 5 8 】

ブロック 6 1 3 において、伝搬器インスタンスは、要求元の指標を記録し、伝搬器インスタンス自体を要求元として示す。伝搬器インスタンスはそれ自体を要求元として示して、ファイルシステムに、応答を伝搬器インスタンスへ返させる。このことは、伝搬器が既存のファイルシステムとインターフェースをとることを容易にする。伝搬器インスタンスは、少なくとも伝搬器インスタンスの観点から実際の要求元を記録し、したがって、下層のファイルシステムからの応答は実際の要求元を示すように更新されることができる。同期関係が半同期である場合、制御はブロック 6 0 9 に流れる。関係が完全同期である場合、制御はブロック 6 1 5 に流れる。

## 【 0 0 5 9 】

ブロック 6 1 5 において、伝搬器インスタンスは、データを記録して飛行中要求を追跡

50

する。「飛行中」要求といわれるが、要求はまだ飛行中でなく、なぜならば、上記データの記録の後まで伝搬器インスタンスが要求を渡さないからである。伝搬器インスタンスは、少なくとも要求の指標と、プライマリエンドポイントと、セカンダリエンドポイントと、応答がプライマリエンドポイント又はセカンダリエンドポイントのいずれかについて受信されたかの指標とを記録する。伝搬器インスタンスは、要求のメタデータから決定される要求の識別子を記録することができる。伝搬器インスタンスは、プライマリエンドポイント識別子及びセカンダリエンドポイント識別子を用いて識別子を生成することができる。伝搬器インスタンスは、このデータを記録して、変更がプライマリエンドポイント及びセカンダリエンドポイントの双方において成功裏に実行されたときを判断する。双方のエンドポイントにおける成功の後、応答が実際の要求元に伝えられることができる。要求を追跡するデータを記録した後、伝搬器インスタンスは、ブロック617及び619を同時に又は順々に実行することができる。順々の場合、伝搬器インスタンスは上記ブロックのうちいずれかを順繰りに実行することができる。

10

**【0060】**

ブロック617において、伝搬器インスタンスは、下層のファイルシステムに要求を供給する。

**【0061】**

ブロック619において、伝搬器インスタンスは、セカンダリエンドポイント上で実行するための要求を示す。例えば、伝搬器インスタンスは、要求及び同期関係情報を、セカンダリエンドポイントに関連付けられたノードに変更を通信するモジュールに渡す。

20

**【0062】**

図7は、論理記憶オブジェクト粒度半同期動作についての例示的な動作のフローチャートを表す。図7は、図6のブロック609からの続きである。

**【0063】**

ブロック701において、伝搬器インスタンスが、データを記録して飛行中要求を追跡する。ブロック701はブロック615と同様の言語において表現されるが、ブロック701の例示的な動作は、セカンダリエンドポイントに関連付けられたノードに送信される要求の状態を追跡しない。セカンダリエンドポイントノードに送信される要求の状態は、半同期関係に関しては伝搬器インスタンスによっては追跡されず、なぜならば、追跡はシーケンサによって扱われるからである。半同期関係において、伝搬器インスタンスは、プライマリエンドポイントについて下層のファイルシステムに送信される要求を追跡することを回避し、下層のファイルシステムの管理メカニズムに依存することができる。この場合、伝搬器インスタンスは、下層のファイルシステムからの成功又は失敗の応答を渡すことができる。伝搬器インスタンスは、応答の中に実際の要求元のアイデンティティを単に復元することになる。

30

**【0064】**

ブロック703において、伝搬器インスタンスは、要求をファイルシステムに供給する。

**【0065】**

図6のブロック609において、シーケンサは、要求及び同期関係情報を渡されている。ブロック705において、シーケンサは、オープン変更セットログの境界が到達されたかを決定する。例えば、境界は、構成されたRPOの断片として定義されてもよい。一例として、シーケンサは、12秒RPOに基づいて、4秒境界上で変更セットログを管理する。変更セットログがオープンされるとき、この変更セットログはシステム時間を用いてスタンプされることができる。境界が到達されるたび、変更セットログはクローズされ、処理のためにサブミットされる。12秒RPOに対する4秒境界の上記例において、合計3つの変更セットログがRPO時間期間に及ぶ。変更セットログのうち1つがオープンであることになり、その他の2つはクローズされることになる。変更セット境界が到達された場合、制御はブロック709に流れる。変更セット境界が到達されなかった場合、制御はブロック707に流れる。

40

50

## 【 0 0 6 6 】

ブロック 7 0 7 において、シーケンサは、オープン変更セットログの中の要求を示す。シーケンサは、要求の識別子、要求のタイプ、及び要求に対する参照を記録することができる。シーケンサは、要求に対する参照を記録することができる。

## 【 0 0 6 7 】

ブロック 7 0 9 において、変更セットログはクローズされ、なぜならば、ブロック 7 0 5 において決定されたとおり境界が到達されたからである。例えば、シーケンサは、オープン変更セットログポインタと 1 つ以上のクローズされた変更セットログポインタとを維持することができる。境界が到達されるとき、シーケンサは、ポインタを更新して、ログのオープン及びクローズを反映することができる。シーケンサは、さらに、ログが失敗の変更セットについてである場合でさえ、クローズされたログをバッファの中に維持することができる。このことは、さらなるメモリを消費する可能性があり、なぜならば、ログが、オープン及びクローズ済みの双方、依然として処理中である変更セットに対して十分なメモリ空間に制約されないからである。しかし、さらなるメモリが、失敗した変更セットの調査を可能にし、あるいは、失敗した変更セットのより速い再試行を容易にすることができる。

10

## 【 0 0 6 8 】

ブロック 7 1 1 において、シーケンサは、クローズされたログを処理のために渡す。例えば、シーケンサは、新しい変更セットログの維持を継続すると同時に、クローズされたログを処理するスレッド又はプロセスを立ち上げることができる。クローズされたログを処理することには、変更要求間における依存関係を決定して変更セットの中の要求間における予期されるシーケンスを満足することが含まれる。

20

## 【 0 0 6 9 】

ブロック 7 1 3 において、シーケンサは、異なる変更セットログをオープンし、初期化する。ログをオープンするために、シーケンサは、異なるメモリ空間を割り振り、あるいは、(成功裏に又は不成功に)完了したクローズされたログのメモリ空間にアクセスすることができる。シーケンサは、オープン変更セットログを初期タイムスタンプを用いて初期化する。シーケンサは、さらに、いかなるデータも上書きして変更セットログをクリアし、あるいは、オープンされる前に別のプロセスがログのクリアを扱うことを可能にすることができる。

30

## 【 0 0 7 0 】

ブロック 7 1 5 において、シーケンサは、初期化されたオープン変更セットログの中の要求を示す。

## 【 0 0 7 1 】

図 8 は、クローズされた変更セットログを処理する例示的な動作のフローチャートを表す。上記処理は、変更セットの中の要求を解析し、任意の順序付けを決定して、正しい及び一貫したデータのビューを維持する。論理記憶オブジェクト粒度が、変更セットをアトミックな仕方を実装し、指定された R P O に準拠するように、変更セットを維持する。図 7 の例示的な動作は、スポンされたスレッド又はプロセスを、クローズされた変更セットログを処理するものとして説明したが、図 8 は、シーケンサを、クローズされた変更

40

## 【 0 0 7 2 】

ブロック 8 0 1 において、シーケンサが、クローズされた変更セットログの中の変更要求間における任意の依存関係を決定し、該依存関係に従ってシーケンシングを示す。シーケンサは、変更要求の的であるエンドポイントの領域を示すデータを維持する。例えば、シーケンサは、ファイルについての領域のビットマップを維持することができる。ビットマップの第 1 の次元は、ファイルシステム及び / 又は記憶プロトコルに依存して、x バイトのブロックを表現することができる。ビットマップの別の次元が、変更要求の各々を表現することができる。上記ビットマップを用いて、シーケンサは、変更要求がいつ重なるかを決定することができる。変更要求が重なる場合、シーケンサは、変更要求が互いに依

50

存すると決定し、そのシーケンスを保存して上記依存関係を満足する。シーケンサは、論理記憶オブジェクトを対象にする任意の変更要求が、該論理記憶オブジェクトを対象にする介在の読み出し要求を有するかをさらに決定する。その場合、シーケンサは、依存関係が存在すると決定し、周りの更新要求のシーケンス又は順序を保存する。シーケンサは、シーケンシング情報を変更要求の各々のメタデータに書き込む。例えば、シーケンサは、シーケンシング情報を変更要求のヘッダに書き込む。

【 0 0 7 3 】

ブロック 8 0 3 において、シーケンサは、変更セットの各変更要求の中の変更セット情報を示す。シーケンシング情報と同様に、シーケンサは、変更要求の各々のメタデータの中の変更セット情報を示す。シーケンシング情報には、変更セットの識別子と変更セット内の変更要求の数とが含まれる。このことは、セカンダリエンドポイントに関連付けられたノードにおけるセカンダリライタが、セカンダリライタが変更セットのすべての変更要求をいつ受信したかを決定するのを助ける。シーケンサは、変更要求の各々のメタデータの中の変更セットの開始時間をさらに示すことができる。このことは、セカンダリライタが、RPO制約がいつ違反されたかを決定するのを助けることができる。

10

【 0 0 7 4 】

ブロック 8 0 5 において、シーケンサは、冗長な変更要求を消去する。シーケンサは、変更要求が後の変更要求によってさらに対象にされる同じプライマリエンドポイント及び同じ領域又はブロックを対象にする場合、変更要求が冗長であると決定する。換言すると、シーケンサは、変更セットを超えて存続するのでない変更を行う変更要求を決定する。

20

【 0 0 7 5 】

ブロック 8 0 7 において、シーケンサは、各変更要求をシーケンシング及び変更セットの指標と共にクラスタ同期エンジンに供給する。シーケンサは、シーケンシング及び変更セット情報の指標を用いて修正された変更要求に対する参照を渡すことができる。こうした変更をセカンダリエンドポイントノードに通信することを担うモジュールが、渡された参照を介して実際のデータを取得することができる。

【 0 0 7 6 】

図 9 は、同期関係におけるプライマリエンドポイントのファイルシステムからの応答を扱う例示的な動作のフローチャートを表す。これら例示的な動作は、伝搬器インスタンスにより実行される場合として説明される。

30

【 0 0 7 7 】

ブロック 9 0 1 において、伝搬器インスタンスが、プライマリエンドポイントのファイルシステムから応答を受信する。伝搬器インスタンスは、変更要求をファイルシステムに事前に渡している。変更要求は、プライマリエンドポイント（すなわち、ファイルハンドル及びファイルブロック番号などのファイル場所情報を有する論理記憶オブジェクト）と、要求のソースとしての伝搬器インスタンスとを示していた。今度は、ファイルシステムが、変更要求をサービス提供した（又は、サービス提供しようとして試みた）後、応答を提供する。応答は、成功又は失敗のいずれかを示すことになる。

【 0 0 7 8 】

ブロック 9 0 3 において、伝搬器インスタンスは、応答が成功又は失敗を示すかを決定する。応答が成功を示す場合、制御はブロック 9 0 9 に流れる。応答が失敗を示す場合、制御はブロック 9 0 5 に流れる。

40

【 0 0 7 9 】

失敗の場合、伝搬器インスタンスは、ブロック 9 0 5 においてセカンダリエンドポイントへの対応する変更の中止（abort）を開始する。プライマリエンドポイントが完全同期関係にあるか又は半同期関係にあるかにかかわらず、セカンダリエンドポイントに対する変更は、プライマリエンドポイントとセカンダリエンドポイントとの間の同期外れ状態（out of sync state）を回避するために、成功裏に完了すべきでない。完全同期関係について、伝搬器インスタンスは、要求を同期エンジンにサブミットして、セカンダリエンドポイントノードに通信される変更要求を中止する。同期エンジンは、セカンダリエンド

50



ポイントに対する変更を中止する動作を実行し、エンドポイント間の同期を保存することになる。半同期関係について、伝搬器インスタンスは、変更セットを中止する。変更セットを中止することには、プライマリエンドポイントノードにおける変更セットログを失敗又は中止としてマーク付けすることと、同期エンジンがセカンダリエンドポイントノードに変更セットの要求を失敗し又は中止するように要求することを要求することを含むことができる。

**【0080】**

ブロック907において、伝搬器インスタンスは、変更要求が失敗したと要求元が通知されることができることを示す。伝搬器インスタンスは、例えば、プライマリエンドポイントのファイルシステムからの失敗応答を変更して実際の要求元を示し、変更された応答をネットワークモジュールに渡すことができる。ネットワークモジュールは、それから、失敗を実際の要求元に通信することができる。

10

**【0081】**

変更要求が成功した場合、伝搬器インスタンスは、ブロック909において追跡データを更新して成功を示す。伝搬器インスタンスは、追跡データを更新して、要求がプライマリエンドポイントにおいて完了したことを示す。

**【0082】**

ブロック911において、伝搬器インスタンスは、追跡データを用いて、セカンダリエンドポイントに対する変更が完了したかを決定する。そうでない場合、制御はブロック913に流れる。セカンダリエンドポイントに対する変更が成功裏に完了した場合、制御はブロック921に流れる。

20

**【0083】**

ブロック913において、伝搬器インスタンスは、タイムアウトが到達されたかを決定する。タイムアウトは構成されることができる。上記タイムアウトは、タイムアウトが満了する前に応答が受信されるべきであることを仮定する。そうでなければ、要求か又は要求に対するセカンダリエンドポイントノードからの応答が、失われたとみなされることができる。タイムアウトが到達された場合、制御はブロック917に流れる。タイムアウトが到達されていない場合、伝搬器インスタンスはブロック915において、定義された待ち期間の間待つ。制御はブロック915から戻ってブロック911に流れる。

**【0084】**

30

ブロック917及び919は、タイムアウトシナリオにおける動作を表す。ブロック917において、伝搬器インスタンスは、セカンダリエンドポイントがプライマリエンドポイントとの同期から外れていることを示す。ブロック919において、エンドポイント間における同期外れ状態は、構成されたとおりに処理される。例えば、再試行が許容される場合、同期外れ状態が再試行を引き起こしてもよい。同期外れ状態が、管理モジュールに対する通知をトリガしてもよい。

**【0085】**

セカンダリエンドポイントにおける更新要求が成功裏に完了するとき、伝搬器インスタンスは、ブロック921において、実際の要求元が変更要求の成功裏の完了を通知されることができることを示す。伝搬器インスタンスは、応答を、実際の要求元に通信するネットワークモジュールに供給する。

40

**【0086】**

ブロック923において、追跡データがクリアされる。伝搬器インスタンスは、このデータをクリアし、あるいは、ガーベッジコレクションスレッドによってクリアされるように追跡データをマーク付けすることができる。

**【0087】**

図9は、プライマリエンドポイントに対する変更の応答を扱う例示的な動作を表すが、図10は、クラスタベースの同期エンジンが伝搬器とカウンターパートの同期エンジンからの要求を処理する例示的な動作のフローチャートを表す。図10の説明は、アクタを同期エンジンという。

50

## 【 0 0 8 8 】

ブロック 1 0 0 1 において、同期エンジンが、セカンダリエンドポイント上で実行されるべき変更要求の指標を受信する。変更要求は、参照で又はリテラルで同期エンジンに渡されることができる。変更要求は、完全同期関係のための変更セットのメンバ、又はスタンドアロンの変更要求であり得る。同期エンジンは、変更要求のメタデータの中、又は変更要求に関連付けられた別個の構造の中に、セカンダリエンドポイントの指標を受信することができる。

## 【 0 0 8 9 】

ブロック 1 0 0 3 において、同期エンジンは、セカンダリエンドポイントに関連付けられたクラスタノードを決定する。同期エンジンは、クラスタにわたり維持されるデータにアクセスする。このデータは、エンドポイント及びノードのためのディレクトリ (directory) として使用されることができる。データは、いずれのノードがいずれの論理記憶オブジェクトに関連付けられる (すなわち、いずれの論理記憶オブジェクトをホストし及び/又はいずれの論理記憶オブジェクトに対するアクセスを管理する) かを示す。上記データは、データベースとして実装されることができる。同期エンジンは、データをセカンダリエンドポイントのアイデンティティを用いて読み出し、該アイデンティティは、論理記憶オブジェクト識別子である。

10

## 【 0 0 9 0 】

ブロック 1 0 0 5 において、同期エンジンは、通信セッションがセカンダリエンドポイントノードにおける同期エンジンとの間ですでに確立されているかを決定する。同期エンジンは、通信セッションを維持して、各要求のための通信セッションを確立するオーバーヘッドを回避する。しかしながら、このことは必要ではない。同期エンジンは、エンドポイントペアごとにセッション又は接続を確立することができる。セッションがすでに確立されてはいない場合、制御はブロック 1 0 0 7 に流れる。そうでなければ、制御はブロック 1 0 0 9 に流れる。

20

## 【 0 0 9 1 】

ブロック 1 0 0 7 において、同期エンジンは、セカンダリエンドポイントに関連付けられたクラスタノードにおける同期エンジンとの通信セッションを確立する。

## 【 0 0 9 2 】

ブロック 1 0 0 9 において、同期エンジンは、変更要求に従ってセカンダリエンドポイントを対象にする複製要求を作成し、伝搬器インスタンスを複製要求のソースとして示す。同期エンジンは、セカンダリエンドポイントを要求の対象として示す要求を作成する。同期エンジンは、セカンダリエンドポイントに書き込まれるべきデータの指標又はデータを有する要求を作成する。同期エンジンは、クラスタノードの指標と受信した変更要求のメタデータとを有する要求をさらに作成する。

30

## 【 0 0 9 3 】

ブロック 1 0 1 1 において、同期エンジンは、セカンダリエンドポイントに関連付けられたクラスタノードに上記セッションを通じて複製要求を通信する。ブロック 1 0 1 1 からブロック 1 0 1 3 への破線は、複製要求の送信と応答の受信との間の時間の経過を表す。

40

## 【 0 0 9 4 】

ブロック 1 0 1 3 において、同期エンジンは、セカンダリノードに関連付けられたクラスタノードから、複製要求に対する応答を受信する。同期エンジンは、応答から伝搬器インスタンスを決定し、上記応答は、伝搬器インスタンスを要求元として示している。同期エンジンは、ブロック 1 0 1 5 において、複製応答を、応答の中に示される適切な伝搬器インスタンスに渡す。同期エンジンは、伝搬器インスタンス識別子を複製要求識別子 (例えば、エンドポイント識別子に基づいて生成される識別子) に関連付けたデータを維持することによって要求元を決定するように設計されることができる。

## 【 0 0 9 5 】

図 1 1 は、伝搬器インスタンスがセカンダリエンドポイントに対する変更要求への応答

50

を扱う例示的な動作のフローチャートを表す。図 11 は、例示的な動作のアクタとして伝搬器インスタンスを参照して説明される。前に説明されたとおり、伝搬器インスタンスは、セカンダリエンドポイントに対して作られることになる変更要求を同期エンジンに渡し、上記変更要求は、複製要求といわれている。同期エンジンは、こうした変更を、セカンダリエンドポイントに関連付けられたクラスタノードに通信する。

【0096】

ブロック 1101 において、伝搬器インスタンスは、同期エンジンから複製応答を受信する。複製応答は、応答がセカンダリエンドポイント及びプライマリエンドポイントに対応することを示す。

【0097】

ブロック 1103 において、伝搬器インスタンスは、エンドポイント間における同期関係が完全同期又は半同期であったかを決定する。同期関係が完全同期である場合、制御はブロック 1105 に流れる。同期関係が半同期である場合、制御はブロック 1123 に流れる。

【0098】

ブロック 1105 において、伝搬器インスタンスは、セカンダリエンドポイントに対する変更要求が成功であったかを、複製応答に基づいて決定する。成功の場合、制御はブロック 1113 に流れる。そうでなければ、制御はブロック 1107 に流れる。

【0099】

ブロック 1107 において、伝搬器インスタンスは、プライマリエンドポイントに対する要求された変更が成功裏に完了したかを決定する。伝搬器インスタンスは、飛行中追跡データを読み出して、プライマリエンドポイント変更が成功裏に完了したかを決定する。プライマリエンドポイントに対する変更が成功裏に完了しており、セカンダリエンドポイントに対する変更が成功しなかった場合、エンドポイントは同期から外れる。プライマリエンドポイントに対する変更が成功裏に完了した場合、制御はブロック 1121 に流れる。プライマリエンドポイントに対する変更が成功裏に完了しなかった場合、制御はブロック 1109 に流れる。

【0100】

ブロック 1121 において、プライマリエンドポイントに対する変更はロールバックされる。プライマリエンドポイントに対する変更をロールバックすることは、要求元が失敗の応答を与えられることにつながる。要求元は、それから、変更を再度要求することができる。伝搬器は、プライマリエンドポイントに対する変更をロールバックすることに追加で又は代わって、同期外れ状態を示すようにプログラムされることができる。制御はブロック 1121 からブロック 1127 に流れる。

【0101】

ブロック 1109 において、プライマリエンドポイントに対する変更は中止される。まれと見込まれるが、セカンダリエンドポイントに対する変更に関連付けられたノードが、伝搬器インスタンスがプライマリエンドポイントの下層の記憶要素から応答を受信する前に、変更要求をサービス提供する可能性がある。

【0102】

ブロック 1111 において、伝搬器インスタンスは、変更要求が失敗したと実際の要求元が通知されることができることを示す。例えば、伝搬器インスタンスは、プライマリエンドポイントの下層の記憶要素からの応答に基づいて、失敗応答を作成する。伝搬器インスタンスは、傍受器から渡された変更要求から事前に記録された要求元を用いて、失敗応答を作成する。

【0103】

セカンダリエンドポイントに対する変更が完全同期関係において成功した場合、伝搬器インスタンスは、ブロック 1113 において、要求された変更がプライマリエンドポイントにおいて成功裏に完了したかを決定する。伝搬器インスタンスは、飛行中追跡データにアクセスして、プライマリエンドポイント変更がすでに完了したかを決定する。プライマ

10

20

30

40

50

リエンドポイント変更がすでに完了している場合、制御はブロック 1 1 1 7 に流れる。プライマリエンドポイント変更がまだ完了していない場合、制御はブロック 1 1 1 5 に流れる。

【 0 1 0 4 】

ブロック 1 1 1 5 において、伝搬器インスタンスは、飛行中追跡データを更新して、セカンダリエンドポイント変更が完了したことを示す。

【 0 1 0 5 】

ブロック 1 1 1 7 において、伝搬器インスタンスは、変更が双方のエンドポイントにおいて成功したとき、要求が完了したと要求元が通知されることができるとを示す。伝搬器インスタンスは、プライマリエンドポイントの下層の記憶要素からの応答に基づいて、応答を生成する。応答は、要求についての成功裏のサービス提供を示す。伝搬器インスタンスは、さらに、要求元としてのそれ自体の指標を、実際の要求元の指標で置換する。伝搬器インスタンスは、それから、応答を傍受器又は通信モジュールに渡す。

【 0 1 0 6 】

ブロック 1 1 1 9 において、要求のための追跡データがクリアされる。伝搬器インスタンスは追跡データをクリアすることができ、あるいは、ガーベッジコレクションスレッド（又は別のデータメンテナンススレッド）が追跡データをクリアすることができる。

【 0 1 0 7 】

応答が、半同期関係におけるセカンダリエンドポイントについてである場合、制御はブロック 1 1 2 3 に流れている。ブロック 1 1 2 3 において、伝搬器インスタンスは、複製応答がセカンダリエンドポイントに対する変更セットの成功裏の完了を示すかを決定する。そうである場合、制御はブロック 1 1 2 5 に流れる。そうでない場合、制御はブロック 1 1 2 7 に流れる。

【 0 1 0 8 】

ブロック 1 1 2 5 において、伝搬器インスタンスは、エンドポイントが同期から外れていることを示す。伝搬器インスタンスは、クラスタノード間で広められる同期関係データにアクセスすることができる。この同期関係データは、伝搬器インスタンスにより設定されることが可能な単一ビットフィールドを含んで、対応するエンドポイントが同期から外れているか又は同期しているかを示すことができる。伝搬エンジンが、同期外れのエンドポイントを伴う要求を、構成されたとおり処理することになる。例えば、伝搬エンジンは、同期外れであると示されるプライマリエンドポイントを対象にするすべての要求を、同期が上記セカンダリエンドポイント又は代替的セカンダリエンドポイントとの間で復元されるまで、さえぎる（fence）ように構成されることができる。伝搬エンジンは、対象にされたエンドポイントが同期外れとして示されるとき、失敗又はサービス外（out of service）タイプの応答で応答するように構成されることができる。プライマリエンドポイントが成功裏に変わり、変更セットがセカンダリエンドポイントにおいて成功裏に完了した場合、制御はブロック 1 1 2 9 に流れる。そうでなければ、制御はブロック 1 1 2 7 に流れる。

【 0 1 0 9 】

ブロック 1 1 2 9 において、伝搬器インスタンスは、変更セットログをクリアする。要求元がプライマリエンドポイントに対する成功裏の変更についてすでに通知されているので、半同期関係におけるセカンダリエンドポイントに対する成功裏の変更は、要求元に対する通知をトリガしない。変更セットログをクリアすることは、変更セットが成功裏に完了したことを暗に示す。伝搬器インスタンスは、クリア又は除去の前、変更セットログを成功裏に完了されたとしてマーク付けするようにプログラムされることができる。

【 0 1 1 0 】

図 1 2 は、セカンダリライタが複製要求を扱う例示的な動作のフローチャートを表す。前に説明されたとおり、セカンダリライタが同期エンジンから複製要求を受信し、その双方がセカンダリエンドポイントに関連付けられたノード上で稼働している。セカンダリエンドポイントとの同期関係にあるプライマリエンドポイントに関連付けられたノード上で

10

20

30

40

50

稼働する同期エンジンが、セカンダリエンドポイントノードにおける同期エンジンに複製要求を通信している。

【0111】

ブロック1201において、セカンダリエンドポイントに関連付けられたクラスタノードにおける同期エンジンが、複製要求を受信する。複製要求は、プライマリエンドポイント及びセカンダリエンドポイントを示す。プライマリエンドポイント又はプライマリエンドポイントノードにおける伝搬器インスタンスが、複製要求のソースとして示される。複製要求は、同期関係のタイプをさらに示すことができる。

【0112】

ブロック1203において、同期エンジンは、セカンダリライタがプライマリ及びセカンダリエンドポイントについてすでにインスタンス化されているかを決定する。例えば、上記エンドポイントペアのためのセカンダリライタが、変更セットについてのより早い要求に対してインスタンス化されている可能性がある。セカンダリライタがすでにインスタンス化されている場合、制御はブロック1207に流れる。そうでない場合、制御はブロック1205に流れる。

10

【0113】

ブロック1207において、同期エンジンは、複製応答と、別個に通信される何らかがある場合には関連するメタデータとを、セカンダリライタインスタンスに渡す。メタデータは、複製要求の中に示されてもよい。

【0114】

ブロック1205において、同期エンジンは、示されるエンドポイントペアに基づいてセカンダリライタをインスタンス化する。制御は、ブロック1205及び1207のいずれかからブロック1209に流れる。

20

【0115】

ブロック1209において、セカンダリライタインスタンスは、変更要求がステージされる(staged)べきかを決定する。変更セットのための変更要求のステージングは、変更要求又は変更要求の指標を制限まで蓄積することを指す。セカンダリライタは、変更要求のメタデータを読み出して、変更要求が変更セット内にあるかを決定することができる。メタデータは、変更セットを示してもよい。セカンダリライタインスタンスは、さらに、完全同期の代わりに半同期の指標に基づいて、変更要求が変更セットの中にある場合として進めることができる。変更要求がステージされることになる場合、制御はブロック1215に流れる。そうでない場合、制御はブロック1211に流れる。

30

【0116】

ブロック1211において、セカンダリライタインスタンスは、データを記録して複製要求を追跡する。セカンダリライタインスタンスは、追跡データを使用して、要求が下層のファイルシステムに渡されていることを記録する。セカンダリライタインスタンスは、データを記録して複製要求を追跡することに代わって、下層のファイルシステムに依存することができる。

【0117】

ブロック1213において、セカンダリライタインスタンスは、複製要求を下層の記憶要素アクセスモジュールに供給する。伝搬器インスタンスと同様に、セカンダリライタインスタンスは、複製要求を下層のファイルシステムに供給する前に、要求元の指標を記録し、それをセカンダリライタインスタンスの指標で置換することができる。

40

【0118】

セカンダリエンドポイントが半同期関係にある場合、セカンダリライタインスタンスは、ブロック1215において、エンドポイントペアのための複製要求をステージすることについて変更セットログがすでに作成されているかを決定することになる。変更セットログがすでに作成されている場合、制御はブロック1217に流れる。変更セットログがエンドポイントペアについてすでに作成されてはいない場合、制御はブロック1225に流れる。

50

## 【 0 1 1 9 】

ブロック 1 2 2 5 において、セカンダリライタインスタンスは、ステージングログ ( staging log ) ( すなわち、セカンダリエンドポイントノードにおける変更セットログ ) を作成する。セカンダリライタインスタンスは、複製要求を用いてステージングログを初期化する。

## 【 0 1 2 0 】

ブロック 1 2 1 7 において、セカンダリライタインスタンスは、すでに作成されたステージングログの中の複製要求を示す。

## 【 0 1 2 1 】

ブロック 1 2 1 9 において、セカンダリライタインスタンスは、変更セットログが完了であるかを決定する。セカンダリライタインスタンスは、複製要求のうちいずれかについてのメタデータにアクセスして、変更セットの中の複製要求の総数を決定することができる。セカンダリライタインスタンスは、それから、メタデータから決定された上記数を、変更セットログ又はステージングログの中に示される複製要求の数に対して比較することができる。ステージングログが完了でない場合、セカンダリライタインスタンスは、さらなる複製要求が受信されるのを待つ。

10

## 【 0 1 2 2 】

ステージングログが完了であるとの決定に対して、セカンダリライタインスタンスは、ブロック 1 2 2 1 においてステージングログを横断する ( traverses ) 。セカンダリライタインスタンスは、ステージングログの中に示される第 1 のマーク付けされていない複製要求を選択し、この選択された複製要求を下層のファイルシステムに供給する。セカンダリライタインスタンスは、対応する応答が受信されるとき、ステージングログ内の次のマーク付けされていない複製要求に進む。セカンダリライタインスタンスは、ステージングログが横断されるまで上記処理を継続し、このことは、図 1 3 においてより詳細に説明される。

20

## 【 0 1 2 3 】

図 1 3 は、セカンダリライタインスタンスが下層のファイルシステムからの応答を扱う例示的な動作のフローチャートを表す。図 1 3 は、変更セットの中の要求の応答を扱う例示的な動作を単に表している。換言すると、図 1 3 は、半同期関係におけるセカンダリエンドポイントの例示的な動作を単に表している。完全同期関係のための応答を扱うとき、セカンダリライタインスタンスは、応答を同期エンジンに渡す。セカンダリライタインスタンスは、最初、元の要求元の指標を復元することになる。

30

## 【 0 1 2 4 】

ブロック 1 3 0 1 において、セカンダリライタインスタンスが、下層のファイルシステムから複製要求に対する応答を受信する。

## 【 0 1 2 5 】

ブロック 1 3 0 3 において、セカンダリライタインスタンスは、応答がセカンダリエンドポイントに対する成功裏の変更を示すかを決定する。応答がセカンダリエンドポイントに対する成功裏の変更を示す場合、制御はブロック 1 3 0 4 に流れる。応答がセカンダリエンドポイントに対する失敗した変更を示す場合、制御はブロック 1 3 0 5 に流れる。

40

## 【 0 1 2 6 】

ブロック 1 3 0 5 において、セカンダリライタインスタンスは、再試行が構成されているかを決定する。セカンダリライタインスタンスは、RPO 適合構成に依存して、変更セット内の要求を再試行するように構成されることができる。例えば、セカンダリライタインスタンスは、所定量の時間が RPO 時間期間内で依然として残っている場合、変更セットを再試行するように構成されることができる。再試行が構成され、許容される場合、制御はブロック 1 3 0 7 に流れる。そうでなければ、制御はブロック 1 3 1 1 に流れる。

## 【 0 1 2 7 】

ブロック 1 3 0 7 において、セカンダリライタインスタンスは、再試行カウンタを更新する。とり得るリソースの無駄を回避するために、再試行は、構成された回数に制限され

50

る。

【0128】

ブロック1309において、セカンダリライタインスタンスは、要求を下層のファイルシステムに再度供給する。

【0129】

再試行が構成されていなかったか、あるいは許容されていなかった場合、セカンダリライタインスタンスは、ブロック1311において、データを記録して、変更セットが失敗したことを示す。セカンダリライタインスタンスは、失敗の指標を、ステージングログのメタデータに書き込むことができる。失敗指標は、失敗がプライマリエンドポイントノードに逆向きに通信されることができないか又は通信されない場合、変更セットについての上記失敗状態を保存するのに役立つ可能性がある。

10

【0130】

ブロック1313において、セカンダリライタインスタンスは、変更セットが失敗したとの通知を生成する。セカンダリライタインスタンスは、変更セットを識別する応答と、失敗の指標とを生成することができる。失敗通知は、それから、同期エンジンを介して要求ノード（すなわち、プライマリエンドポイントノード）に供給される。ブロック1313からブロック1315への点線は、時間の経過を示す。より後の時間に、セカンダリライタインスタンスは、変更セットログを破棄のためにマーク付けすることができる。セカンダリライタは、失敗の指標を用いて変更セットログを破棄するようにプログラムされることができる。

20

【0131】

複製要求応答が成功を示した場合、セカンダリライタインスタンスは、ブロック1304において、対応する変更が完了したかを決定する。変更セット全体が完了している場合、制御はブロック1317に流れる。変更セット全体が完了してはいない場合、制御はブロック1319に流れる。

【0132】

ブロック1317において、セカンダリライタインスタンスは、変更セットが成功裏に完了したとの通知を生成する。セカンダリライタインスタンスは、変更セットを識別する応答と、成功の指標とを生成することができる。成功通知は、それから、同期エンジンを介して要求ノード（すなわち、プライマリエンドポイントノード）に供給される。セカンダリライタは、変更セット成功通知を生成することに代わって、変更セット内の複製要求のうち1つに対する応答を渡して戻すようにプログラムされることができる。セカンダリライタは、変更セットのシーケンシング情報に従って、最後の要求に対する応答を要求ノードに返すことができる。変更セットの最後の変更要求に対するこの応答が、プライマリエンドポイントノードにおける伝搬器インスタンスに対して変更セット全体の成功通知として動作することができる。

30

【0133】

ブロック1319において、セカンダリライタインスタンスは、ステージングログの中の個別の要求を成功裏に完了したとしてマーク付けする。

【0134】

ブロック1321において、セカンダリライタインスタンスは、ステージングログの横断を継続する。セカンダリライタインスタンスは、ステージングログ内の要求について示される順序付けに従って、ステージングログ内の次のマーク付けされていない要求を選択する。セカンダリライタインスタンスは、この選択された要求を下層のファイルシステムに供給する。

40

【0135】

例示からのバリエーション

フローチャートは、例示の理解の助けとなるように提供されており、請求項の範囲を限定するのに使用されるべきではない。フローチャートは、開示の態様間で変動し得る例示的な動作を表す。さらなる動作が実行されてもよく、より少ない動作が実行されてもよく

50

、動作が並列に実行されてもよく、動作が異なる順序で実行されてもよい。例えば、傍受器が、要求を伝搬器に単に渡してもよい。伝搬器は、変更要求の対象が同期関係にあるかを決定するようにプログラムされることができる。バリエーションの別の例として、ブロック615が、同期構成にかかわらず実行されることができる。追跡データの欠如が半同期関係を暗に示すこと及び応答が逆向きに要求元に供給されることに代わって、各要求の追跡データが維持されることができる。フローチャートの動作のいくつかは、スレッド又はプロセスがエンドポイントペアについてすでにインスタンス化されているかを決定することを説明した。存続スレッドを使用しないアーキテクチャが、設計されることができる。代わりに、状態データがエンドポイントペアごとに記憶される。この状態データは、対応する（1つ以上の）要求が完了し又は失敗した後にクリアされるまで存続する。このことは、リソースを消費する待ち状態のスレッドを回避する。別の例として、伝搬器インスタンスが、シーケンサに代わって又はシーケンサに対して追加で、セカンダリエンドポイントに対する要求の状態を追跡することができる。変更セットの状態を追跡する個別のアクタにかかわらず、変更セット内の個々の要求の状態は、変更がプライマリエンドポイントにおいて実行された後に要求元に応答することを妨げない。ゆえに、変更セット内の個々の要求の状態は、追跡される必要がない。図7において、ブロック711は、クローズされた変更セットログをプロセスのスポンされたスレッドに渡すことを説明している。オープンログを処理するアーキテクチャがプログラムされ、あるいは設計されることができる。アーキテクチャは、要求が追加されるたび変更セットログを処理し、シーケンシング情報を更新し、冗長な変更を消去するなどすることができる。変更セットログがクローズされるとき、該ログはすでに順序付けられており、セカンダリエンドポイントに関連付けられたノードに通信する準備ができています。

10

20

30

40

50

**【0136】**

失敗の通信は、フローチャート内に表される仕方とは異なる仕方で通信されることができる。例えば、失敗指標（例えば、ブロック1311）は記録されなくてもよく、なぜならば、失敗通知が生成されるからである。失敗した変更セットについて、セカンダリライタインスタンスが、変更セット内の要求のうち1つについての失敗の応答を渡して戻すことができる。要求ノードにおける伝搬器インスタンスは、いずれのクローズされた変更セットログが上記失敗応答に対応するかを決定し、その変更セットログを失敗したとしてマーク付けすることができる。

**【0137】**

さらに、表されていないさらなる動作が実行されることができる。例えば、変更セットログを監視する監視スレッドがスポンされることことができる。監視スレッドは、定義されたRPOに対して、変更セットの寿命を評価することができる。アクティブな変更セットログ又は処理中の変更セットログは、プライマリエンドポイントノード又はセカンダリエンドポイントノードのいずれかからの応答を依然として待つクローズされた変更セットログである。監視スレッドは、変更セット開始時間を評価して、RPO時間が経過したかを決定する。そうである場合、監視スレッドは、変更セットを失敗したとしてマーク付けするようにシーケンサスレッドに促し、あるいは、それ自体で変更セットを失敗したとしてマーク付けすることができる。

**【0138】**

本説明は、同期関係についてペアにされている個々の論理記憶オブジェクトを参照するが、同期関係の「エンドポイント」は、論理記憶オブジェクトのグループとすることができる。例えば、ファイルのグループ又はLUNのグループが、論理記憶オブジェクトの別のグループとの同期関係にあることができる。ノードは、さらなるデータを維持して、グループのメンバである論理記憶オブジェクトに対するグループ識別子を解決することができる。

**【0139】**

当業者には十分理解されるであろうとおり、本開示の態様は、システム、方法、又はコンピュータプログラム製品として実施されることができる。したがって、本開示の態様は



、ハードウェア態様、ソフトウェア態様（ファームウェア、常駐ソフトウェア、マイクロコード等を含む）、又はソフトウェア態様とハードウェア態様とを組み合わせた態様の形式をとることができ、これら態様はすべて一般に「回路」、「モジュール」、又は「システム」といわれることがある。さらに、本開示の態様は、コンピュータ読取可能プログラムコードを具現化させた1つ以上のコンピュータ読取可能媒体において具現化されたコンピュータプログラム製品の形式をとってもよい。

【0140】

1つ以上のコンピュータ読取可能媒体の任意の組み合わせが利用されてもよい。コンピュータ読取可能媒体は、コンピュータ読取可能信号媒体又はコンピュータ読取可能記憶媒体であり得る。コンピュータ読取可能記憶媒体は、例えば、これらに限られないが、電子的、磁氣的、光学的、電磁的、赤外線、又は半導体の、システム、装置、若しくはデバイス、又は前述のものの任意の適切な組み合わせであり得る。コンピュータ読取可能記憶媒体のさらなる具体的な例（包括的でないリスト）には、下記が含まれる：1つ以上のワイヤを有する電気接続、ポータブルコンピュータディスク、ハードディスク、ランダムアクセスメモリ（RAM）、読取専用メモリ（ROM）、消去可能プログラマブル読取専用メモリ（EPROM又はフラッシュメモリ）、光ファイバ、ポータブルコンパクトディスク読取専用メモリ（CD-ROM）、光学記憶装置、磁気記憶装置、又は前述のものの任意の適切な組み合わせ。本文献の文脈において、コンピュータ読取可能記憶媒体は、命令実行システム、装置、若しくはデバイスによってか又はこれらに関連して使用されるプログラムを含み又は記憶することができる任意の有形媒体であり得る。

10

20

【0141】

コンピュータ読取可能信号媒体は、コンピュータ読取可能プログラムコードを例えばベースバンドにおいて又は搬送波の一部として具現化された、伝搬されるデータ信号を含み得る。こうした伝搬信号は、様々な形式のうち任意のものをとることができ、これらに限られないが、電子磁気信号、光信号、赤外線信号、又はこれらのうち任意の適切な組み合わせが含まれる。コンピュータ読取可能信号媒体は、コンピュータ読取可能記憶媒体でなく、コンピュータによってか又はコンピュータと関連して使用されるプログラムを通信し、伝搬し、又は輸送することができる、任意のコンピュータ読取可能媒体であり得る。コンピュータ読取可能信号媒体上に具現化されたプログラムコードは、任意の適切な媒体を用いて伝達されてよく、これらに限られないが、無線、有線、光ファイバケーブル、RF等、又は前述のものの任意の適切な組み合わせが含まれる。

30

【0142】

本開示の態様の動作を実行するコンピュータプログラムコードは、1つ以上のプログラミング言語の任意の組み合わせにおいて書かれてよく、プログラミング言語には、オブジェクト指向プログラミング言語、例えばJava（登録商標）プログラミング言語、C++など、ダイナミックプログラミング言語、例えばPythonなど、スクリプト言語、例えばPerlプログラミング言語又はPowerShell（登録商標）スクリプト言語など、会話手続型プログラミング言語、例えば“C”プログラミング言語又は同様のプログラミング言語などが含まれる。プログラムコードは、全体的にスタンドアロンコンピュータ上で実行されてもよく、複数のコンピュータにわたり分散される仕方で行われてもよく、1つのコンピュータ上で実行されると同時に別のコンピュータ上で結果を提供し又は入力を受け入れてもよい。

40

【0143】

本開示の態様は、本開示の態様に従う方法、装置（システム）、及びコンピュータプログラム製品のフローチャート例示及び/又はブロック図を参照して説明される。フローチャート例示及び/又はブロック図の各ブロック、並びにフローチャート例示及び/又はブロック図内の複数ブロックの組み合わせが、コンピュータプログラム命令によって実装できることが理解されるであろう。上記コンピュータプログラム命令は、汎用目的コンピュータ、特別目的コンピュータ、又は他のプログラマブルデータ処理装置のプロセッサに提供されてマシンを生じさせることができ、したがって、コンピュータ又は他のプログラマ

50

ブルデータ処理装置のプロセッサを介して実行される上記命令は、フローチャート及び／又はブロック図の1つ又は複数のブロック内に指定される機能／動作を実施する手段を作り出す。

【0144】

上記コンピュータプログラム命令は、コンピュータ、他のプログラマブルデータ処理装置、又は他のデバイスに特定の仕方で機能するように指示できるコンピュータ記憶媒体内にさらに記憶されてもよく、したがって、コンピュータ読取可能媒体内に記憶された上記命令は、フローチャート及び／又はブロック図の1つ又は複数のブロック内に指定される機能／動作を実施する命令を含む製造品を生じさせる。

【0145】

コンピュータプログラム命令は、さらに、コンピュータ、他のプログラマブルデータ処理装置、又は他のデバイスにロードされて、一連の動作ステップをコンピュータ、他のプログラマブル装置、又は他のデバイス上で実行させて、コンピュータ実施処理を生じさせてもよく、したがって、コンピュータ又は他のプログラマブル装置上で実行される上記命令は、フローチャート及び／又はブロック図の1つ又は複数のブロック内に指定される機能／動作を実施する処理を提供する。

【0146】

図14は、記憶クラスタベースの完全同期及び半同期伝搬エンジンを備えた一例示的なコンピュータシステムを表す。コンピュータシステムは、プロセッサユニット1401（可能性として複数のプロセッサ、複数のコア、複数のノードを含み、かつ／あるいはマルチスレッディングを実施するなどする）を含む。コンピュータシステムは、メモリ1407を含む。メモリ1407は、システムメモリ（例えば、キャッシュ、SRAM、DRAM、ゼロキャパシタRAM、ツイントランジスタRAM、eDRAM、EDORAM、DDR RAM、EEPROM、NRAM、RRAM（登録商標）、SONOS、PRAM等のうち1つ以上）又は上記ですでに説明されたマシン読取可能媒体のとり得る具現化のうち任意の1つ以上であり得る。コンピュータシステムは、バス1403（例えば、PCI、ISA、PCIエクスプレス、HyperTransport（登録商標）バス、InfiniBand（登録商標）バス、NuBus等）と、ネットワークインターフェース1405（例えば、ATMインターフェース、イーサネット（登録商標）インターフェース、フレームリレーインターフェース、SONETインターフェース、無線インターフェース、iSCSI、ファイバチャネル等）とをさらに含む。コンピュータシステムは、記憶クラスタベース粒状（storage cluster based granular）完全同期及び半同期伝搬エンジン1411をさらに含む。記憶クラスタベース粒状完全同期及び半同期伝搬エンジン1411は、上記で説明されたとおり、同期関係におけるエンドポイントを対象にするファイルシステム変更要求に対応する要求及び応答を扱う。こうした機能性のうち任意のものが、ハードウェアにおいて及び／又は処理ユニット1401上で部分的に（又は全体的に）実施されてよい。例えば、上記機能性は、特定用途向け集積回路を用いて、処理ユニット1401内に実装されるロジック内で、周辺デバイス又はカード上のコプロセッサ内でなどで実施されてよい。さらに、具現化には、より少ないコンポーネント、又は図14内に例示されていないさらなるコンポーネント（例えば、ビデオカード、オーディオカード、さらなるネットワークインターフェース、周辺デバイス等）を含んでもよい。処理ユニット1401、（1つ以上の）記憶装置1409、及びネットワークインターフェース1405は、バス1403に結合される。メモリ1407は、バス1403に結合されるものとして例示されているが、処理ユニット1401に結合されてもよい。

【0147】

本開示の態様は、様々な実装及び活用を参照して説明されているが、本開示のこうした態様は例示であり、発明対象事項の範囲はこれらに限定されないことが理解されるであろう。一般に、本明細書に説明されるクラスタノードにわたる論理記憶オブジェクト粒度同期の手法は、任意の1つ又は複数のハードウェアシステムに調和する設備を用いて実装されてよい。多くのバリエーション、修正、追加、及び改善が可能である。

10

20

30

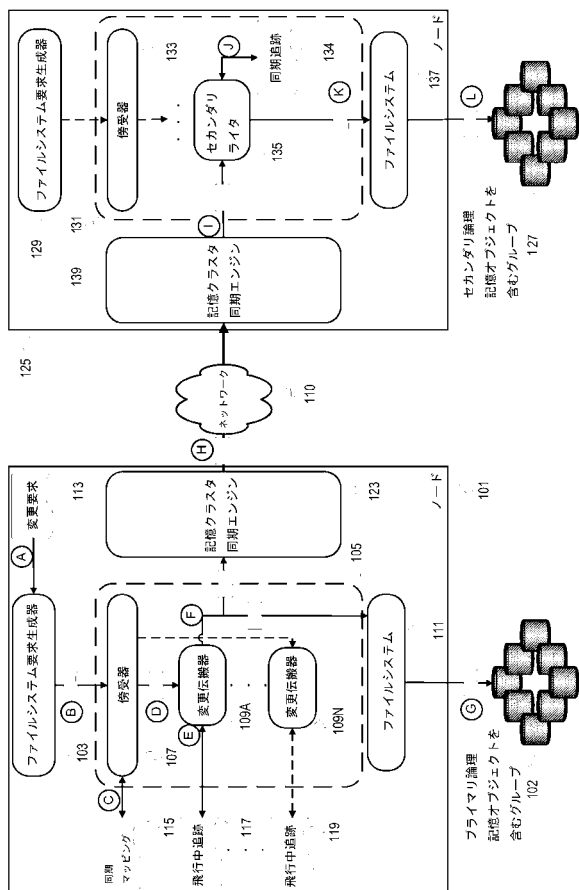
40

50

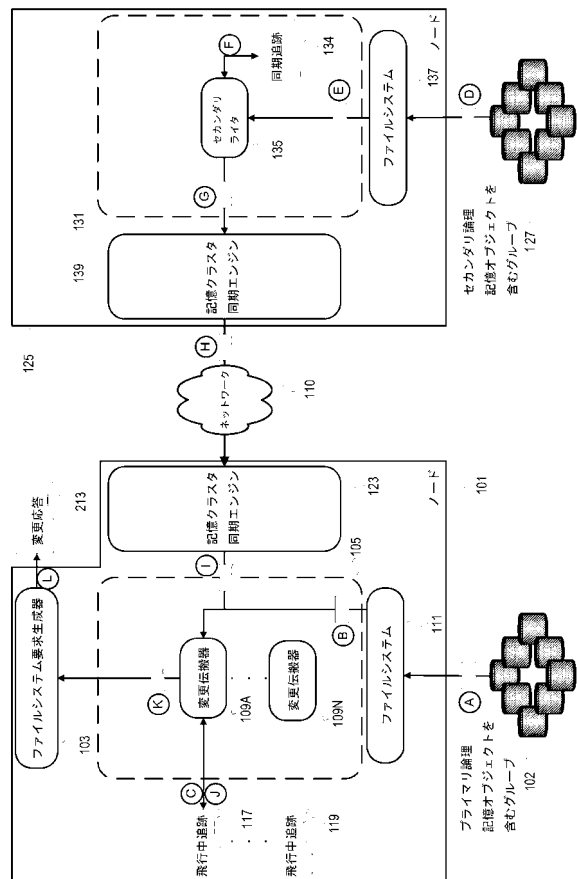
【0148】

本明細書において単一のインスタンスとして説明されるコンポーネント、動作、又は構造について、複数のインスタンスが提供されてもよい。最後、様々なコンポーネント、動作、及びデータストア間の境界はいくらか任意的であり、特定の動作が具体的な例示構成の文脈の中で例示されている。機能性についての他の割り振りが想定され、発明対象事項の範囲内に入り得る。一般に、例示構成の中で別個のコンポーネントとして提示される構造及び機能性は、組み合わせられた構造又はコンポーネントとして実装されてもよい。同様に、単一のコンポーネントとして提示される構造及び機能性が、別個のコンポーネントとして実装されてもよい。上記及び他のパリエーション、修正、追加、及び改善が発明対象事項の範囲内に入り得る。

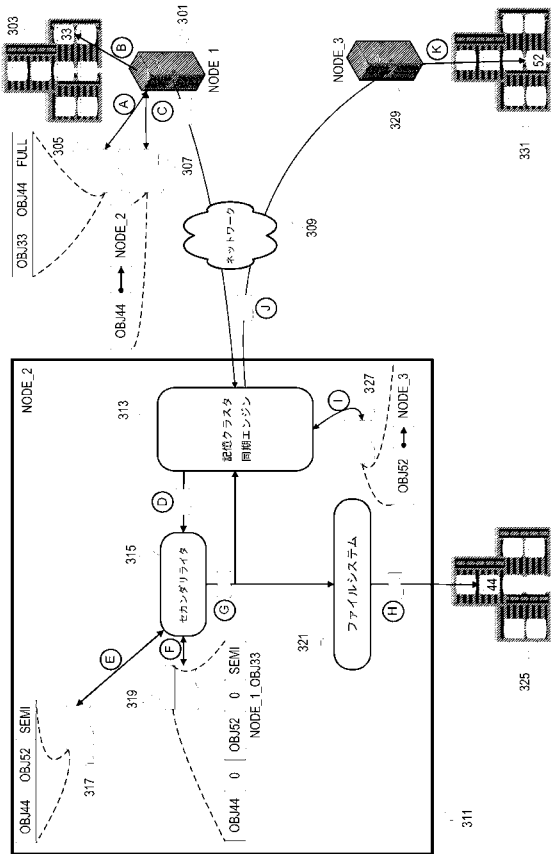
【図1】



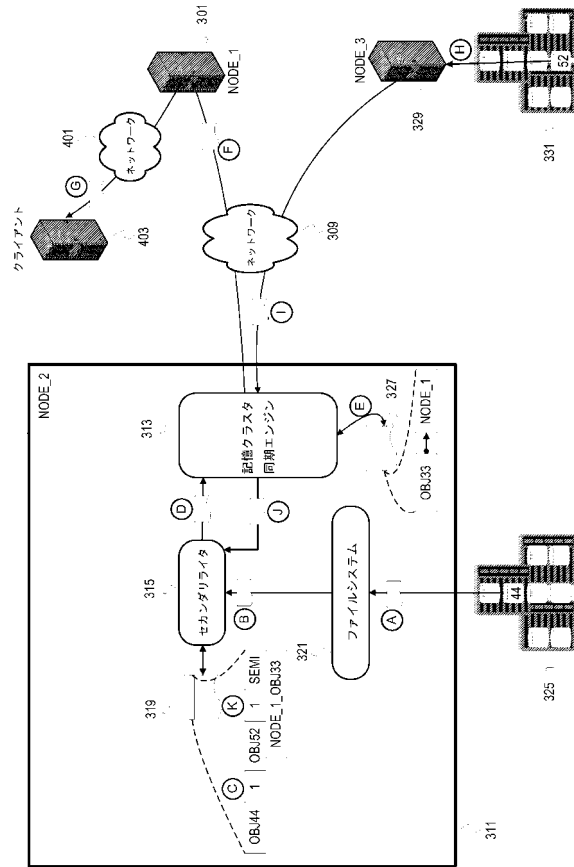
【図2】



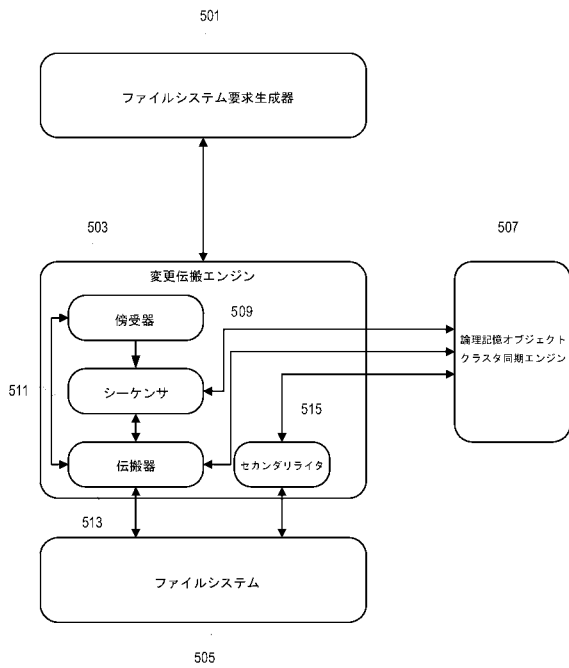
【図3】



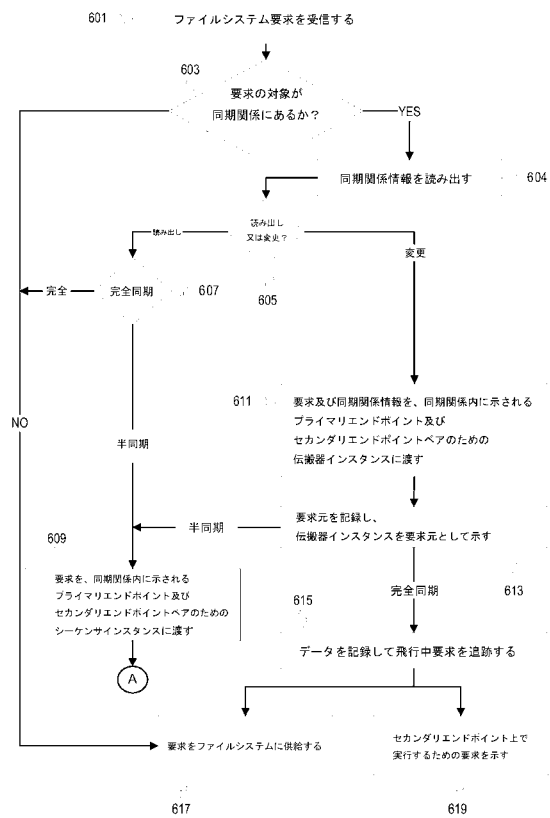
【図4】



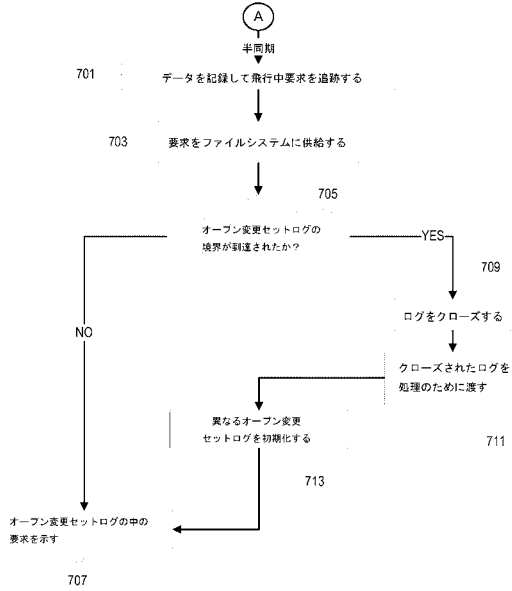
【図5】



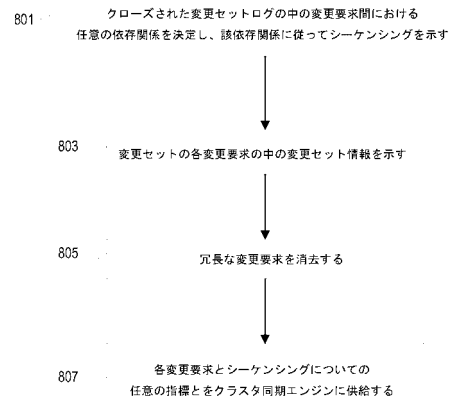
【図6】



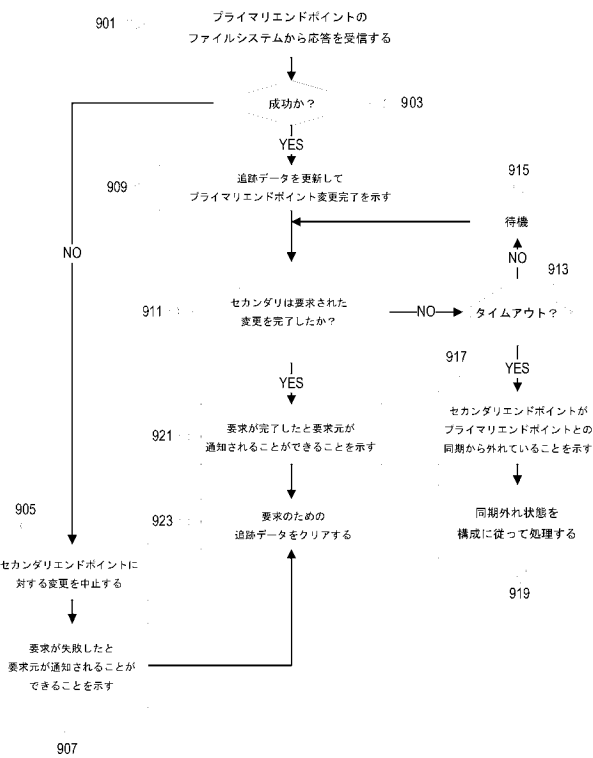
【 図 7 】



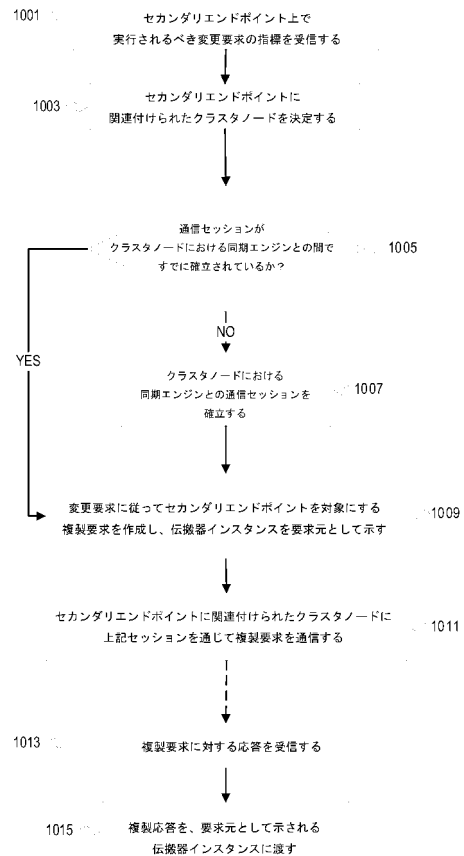
【 図 8 】



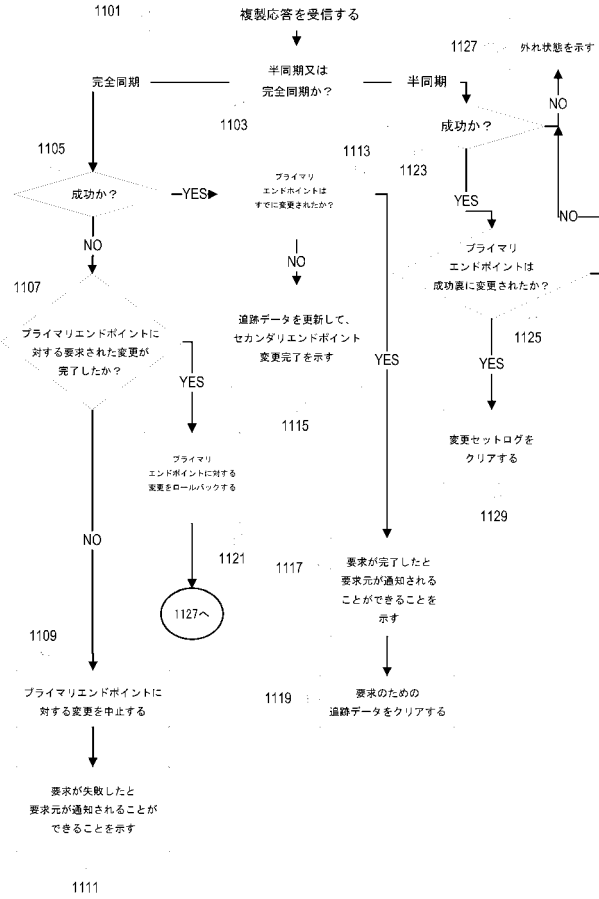
【 図 9 】



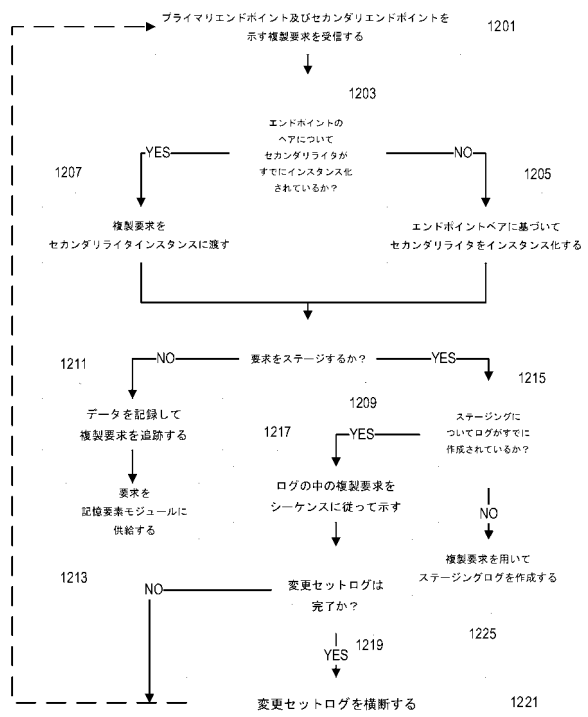
【 図 10 】



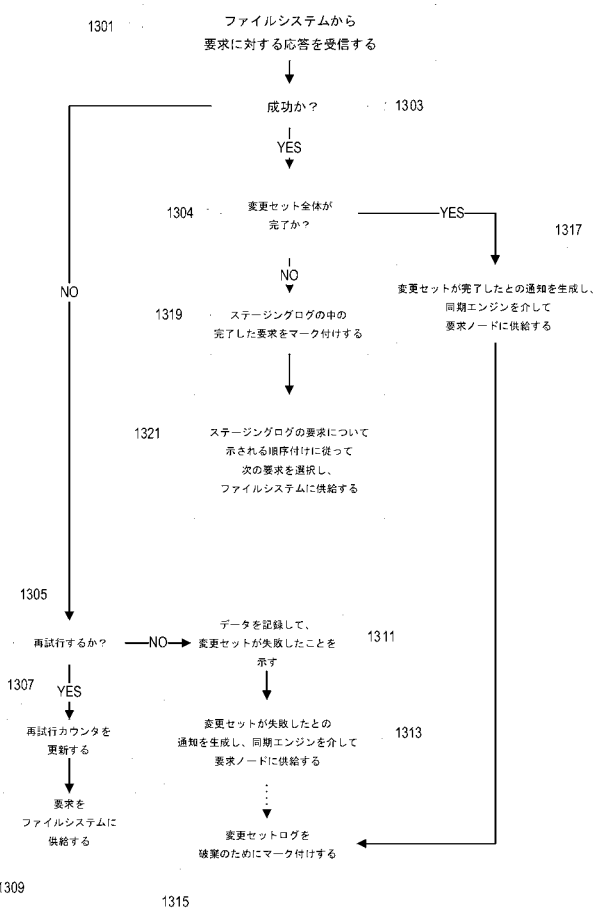
【図 1 1】



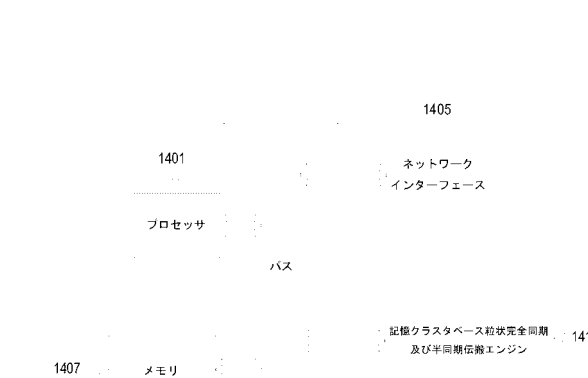
【図 1 2】



【図 1 3】



【図 1 4】



## 【 国際調査報告 】

## INTERNATIONAL SEARCH REPORT

International application No PCT/US2015/043159
---

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> INV. G06F17/30 H04L29/08 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) G06F H04L		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal, WPI Data, COMPENDEX		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 7 395 352 B1 (LAM SAHN [US] ET AL) 1 July 2008 (2008-07-01) column 1, line 51 - line 62 column 3, line 48 - line 60 -----	1-20
X	US 2011/066592 A1 (NEWPORT WILLIAM T [US] ET AL) 17 March 2011 (2011-03-17) page 1, left-hand column, paragraph 0008 - page 1, left-hand column, paragraph 0008 figure 4 page 3, right-hand column, paragraph 0030 ----- -/--	1-20
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents : "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search		Date of mailing of the international search report
20 October 2015		28/10/2015
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer
		Stan, Johann

4

## INTERNATIONAL SEARCH REPORT

International application No

PCT/US2015/043159

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 7 039 661 B1 (RANADE DILIP M [IN]) 2 May 2006 (2006-05-02) column 3, line 5 - line 13 column 7, line 58 - line 64 column 8, line 29 - line 35 column 5, line 18 - line 21 column 5, line 25 - line 45 -----	2,6,10, 14,17
A	US 8 495 250 B2 (ANANTHANARAYANAN RAJAGOPOL [US] ET AL) 23 July 2013 (2013-07-23) column 7, line 4 - line 26 -----	2,6,10, 14,17
A	US 8 667 236 B2 (PHELPS ADAM M [US] ET AL) 4 March 2014 (2014-03-04) column 2, line 57 - column 3, line 5 -----	2,6,10, 14,17
A	US 2005/193041 A1 (BOURBONNAIS SERGE [US] ET AL) 1 September 2005 (2005-09-01) the whole document -----	2,6,10, 14,17
A	US 7 571 391 B2 (ROESSLER ANDREAS [DE]) 4 August 2009 (2009-08-04) the whole document -----	3,11,18
A	"Sun StorageTek (TM) Availability Suite 4.0 Remote Mirror Software Administration Guide",  30 April 2006 (2006-04-30), XP055220220, Retrieved from the Internet: URL:https://docs.oracle.com/cd/E19359-01/8 19-6148-10/819-6148-10.pdf [retrieved on 2015-10-12] page 19 - page 32 -----	1,8,9,16



**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No

PCT/US2015/043159

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 7395352	B1	01-07-2008	NONE
US 2011066592	A1	17-03-2011	JP 5496839 B2 21-05-2014 JP 2011060292 A 24-03-2011 KR 20110029071 A 22-03-2011 US 2011066592 A1 17-03-2011 US 2013041869 A1 14-02-2013
US 7039661	B1	02-05-2006	US 7039661 B1 02-05-2006 US 7606841 B1 20-10-2009
US 8495250	B2	23-07-2013	US 2011145499 A1 16-06-2011 US 2012311065 A1 06-12-2012
US 8667236	B2	04-03-2014	CN 103124961 A 29-05-2013 EP 2622487 A2 07-08-2013 US 2012079222 A1 29-03-2012 WO 2012050904 A2 19-04-2012
US 2005193041	A1	01-09-2005	US 2005193041 A1 01-09-2005 US 2007288537 A1 13-12-2007 US 2008163222 A1 03-07-2008
US 7571391	B2	04-08-2009	NONE

## フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(72)発明者 コートニー , スーザン , エム .

アメリカ合衆国 カリフォルニア州 9 4 0 8 9 サニーベール イースト ジャヴァ ドライヴ  
4 9 5

(72)発明者 ムウ , ユエドーン

アメリカ合衆国 カリフォルニア州 9 4 0 8 9 サニーベール イースト ジャヴァ ドライヴ  
4 9 5

(72)発明者 ラオ , サントシュ

アメリカ合衆国 カリフォルニア州 9 4 0 8 9 サニーベール イースト ジャヴァ ドライヴ  
4 9 5