

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5426380号
(P5426380)

(45) 発行日 平成26年2月26日(2014.2.26)

(24) 登録日 平成25年12月6日(2013.12.6)

(51) Int.Cl.

F I

G 0 6 F 12/00 (2006.01)

G 0 6 F 12/00 5 3 1 D

G 0 6 F 12/00 5 1 7

請求項の数 14 (全 37 頁)

(21) 出願番号	特願2009-525761 (P2009-525761)	(73) 特許権者	506329306
(86) (22) 出願日	平成19年8月22日 (2007.8.22)		アマゾン テクノロジーズ インコーポレ
(65) 公表番号	特表2010-501942 (P2010-501942A)		イテッド
(43) 公表日	平成22年1月21日 (2010.1.21)		アメリカ合衆国 89507 ネバダ州
(86) 国際出願番号	PCT/US2007/076535		レノ ビーオー ボックス 8102
(87) 国際公開番号	W02008/024850	(74) 代理人	100061815
(87) 国際公開日	平成20年2月28日 (2008.2.28)		弁理士 矢野 敏雄
審査請求日	平成22年8月23日 (2010.8.23)	(74) 代理人	100094798
(31) 優先権主張番号	11/508, 129		弁理士 山崎 利臣
(32) 優先日	平成18年8月22日 (2006.8.22)	(74) 代理人	100099483
(33) 優先権主張国	米国 (US)		弁理士 久野 琢也
		(74) 代理人	100112793
			弁理士 高橋 佳大
		(74) 代理人	100128679
			弁理士 星 公弘

最終頁に続く

(54) 【発明の名称】 高可用性データを提供するためのシステム及び方法

(57) 【特許請求の範囲】

【請求項 1】

プロセッサと、
メモリと、

を備えているコンピュータ実施データ記憶システムにおいて、

前記メモリは、データ記憶装置を提供するように前記システムを構成する一つ又は複数のコンピュータ実行可能ロジックコンポーネントを備え、

前記コンピュータ実行可能ロジックコンポーネントは、

リング接続構成に対応する複数のデータセンタに及び前記複数のデータセンタ内の複数のホストに複数のデータセットを記憶するための分担をマッピングするように構成されたマッピングロジックと、

前記複数のデータセンタの第1のサブセット内の前記複数のホストの第1のサブセットに、データセットの第1の複数のコピーを書き込むように構成されたデータセット複製ロジックであって、前記複数のホストの前記第1のサブセットは前記データセットのための優先リストに対応し、前記複数のデータセンタの前記第1のサブセットは前記リング接続構成上における所定の位置により識別され、複数部分キーの関数として識別される前記所定の位置は前記データセットに対応する、データセット複製ロジックと、

前記複数のデータセンタの第2のサブセット内の前記複数のホストの第2のサブセットで、前記データセットの第2の複数のコピーを読み出すことにより前記データセットの一つのコピーの供給の要求に応答するように構成されたデータセット検索ロジックであって

10

20

、前記複数のホストの前記第 2 のサブセットは前記データセットのための優先リストに対応する、データセット検索ロジックと、

前記データセットの一つのコピーを供給するために、前記データセットの前記第 2 の複数のコピー間の因果関係を評価するようにデータセット比較ロジックと、
を備え、

因果関係の前記評価は、

前記データセットの前記第 2 の複数のコピーのうちのコピーであって、前記データセットの前記第 2 の複数のコピーのうちの他のコピーの先祖でないコピーに規定の順序を適用すること、

前記データセットの前記第 2 の複数のコピーのうちのコピーであって、前記データセットの前記第 2 の複数のコピーのうちの他のコピーの先祖でないコピーを併合すること、及び、

前記データセットの前記第 2 の複数のコピーのうちのコピーであって、前記データセットの前記第 2 の複数のコピーのうちの他のコピーの先祖でないコピーの間の矛盾に特定の調整アルゴリズムを適用すること、

のうちの少なくとも一つを含み、

前記データセットの書込のための前記複数のホストの前記第 1 のサブセットと前記データセットの読出のための前記複数のホストの前記第 2 のサブセットとは、別個に決定されることを特徴とするコンピュータ実施データ記憶システム。

【請求項 2】

前記マッピングロジックは、ハッシュ関数に基づいてハッシュ値を生成するように構成されたロジックを含み、個々の前記データセンタが、前記ハッシュ範囲の複数の異なる部分にある前記複数のデータセットのサブセットを記憶する分担を有するように、前記複数のデータセンタのそれぞれは、前記ハッシュ関数のハッシュ範囲内にある複数の位置を有することを特徴とする請求項 1 に記載のシステム。

【請求項 3】

前記ハッシュ値は第 1 のハッシュ値であり、前記ハッシュ関数は第 1 のハッシュ関数であり、前記ハッシュ範囲は第 1 のハッシュ範囲であり、前記マッピングロジックは、第 2 のハッシュ範囲における第 2 のハッシュ関数に基づいて第 2 のハッシュ値を生成するように構成されたロジックを含むことを特徴とする請求項 2 に記載のシステム。

【請求項 4】

個々の前記データセンタが、前記第 2 のハッシュ範囲の複数の異なる部分にある前記複数のデータセットのサブセットを記憶する分担を有するように、前記データセンタのそれぞれは、前記第 2 のハッシュ範囲内にある複数の位置を有することを特徴とする請求項 3 に記載のシステム。

【請求項 5】

データ貸出の終了後に前記データセットの他のコピーを更新するように構成されている貸出ロジックをさらに備えていることを特徴とする請求項 1 に記載のシステム。

【請求項 6】

前記データセット検索ロジックは、前記データセットの前記第 2 の複数のコピーをブリフェッチするように構成されていることを特徴とする請求項 1 に記載のシステム。

【請求項 7】

前記データセット比較ロジックは、前記データセットの前記第 2 の複数のコピーのそれぞれに関連して記憶されたバージョン履歴に基づいて前記因果関係を評価するように構成されており、前記バージョン履歴は、それぞれのハッシュ履歴及びそれぞれのベクトルクロックのうちの少なくとも一方を含むことを特徴とする請求項 1 に記載のシステム。

【請求項 8】

リング接続構成に対応する複数のデータセンタに及び前記複数のデータセンタ内の複数のホストに複数のデータセットを記憶するための分担を、各データセットに対応する複数部分キーのための第 1 及び第 2 のハッシュ関数を用いてマッピングするステップであって

10

20

30

40

50

、前記第 1 のハッシュ関数は、前記複数のデータセンタのうちの選択されたデータセンタに前記複数のデータセットを記憶するための分担をマッピングし且つ前記複数部分キーの第 1 の部分を入力として使用し、前記第 2 のハッシュ関数は、前記選択されたデータセンタ内の前記複数のホストのうちの選択されたホストに前記複数のデータセットを記憶するための分担をマッピングし且つ前記複数部分キーの第 2 の部分を入力として使用する、ステップと、

データセットの複数のコピーのそれぞれについてのバージョン履歴の書込を含めて、前記複数のデータセンタのサブセットに、前記複数のデータセンタの前記サブセットの所定の位置に基づいて前記データセットの前記複数のコピーを記憶するステップと、

前記データセットの前記複数のコピーについてのバージョン履歴に基づいて、前記データセットのコピー間の因果関係を評価するステップと、

を含み、

因果関係の前記評価は、
前記データセットの前記複数のコピーのうちのコピーであって、前記データセットの前記複数のコピーのうちの他のコピーの先祖でないコピーに規定の順序を適用すること、

前記データセットの前記複数のコピーのうちのコピーであって、前記データセットの前記複数のコピーのうちの他のコピーの先祖でないコピーを併合すること、及び、

前記データセットの前記複数のコピーのうちのコピーであって、前記データセットの前記複数のコピーのうちの他のコピーの先祖でないコピーの間の矛盾に特定の調整アルゴリズムを適用すること、

のうちの少なくとも一つを含む、

ことを特徴とする、

各ステップをプロセッサに実行させるためのコンピュータ実施データ処理方法。

【請求項 9】

前記バージョン履歴の書込は、前記複数のデータセンタの前記サブセットの優先リストに従って実行され、前記優先リストは、ハッシュ関数に基づくものであることを特徴とする請求項 8 に記載の方法。

【請求項 10】

ハッシュキー及び前記ハッシュ関数に基づいてハッシュ値を生成するステップをさらに含み、前記ハッシュキーは、前記データセットに関連するものであり、前記ハッシュ関数に入力として適用され、前記ハッシュ関数は、前記複数のデータセットを記憶するための分担をマッピングし、前記複数のデータセンタの前記サブセットは、前記ハッシュ値に基づいて、及び、他のデータセンタが利用不可能であるか否かに基づいて、前記データセットを記憶するために選択されることを特徴とする請求項 9 に記載の方法。

【請求項 11】

個々の前記データセンタが、前記リング接続構成内にある複数の位置に対応する前記ハッシュ範囲の複数の異なる部分にある前記複数のデータセットのサブセットを記憶する分担を有するように、前記データセンタのそれぞれは、前記ハッシュ範囲内にある前記ハッシュ値に基づく所定の位置を有することを特徴とする請求項 10 に記載の方法。

【請求項 12】

前記データセットの前記複数のコピーの前記記憶は、優先リストに従って実行され、前記優先リストは、前記データセットの前記複数のコピーが記憶されるべきデータセンタの順位を規定し、予め割り当てられた値が、書込が実行されることとなる、前記優先リスト内にあるデータセンタの数を決定し、

前記方法は、前記データセットの前記複数のコピーのうちの一つを第 1 のデータセンタから第 2 のデータセンタへ、前記第 2 のデータセンタが利用可能になった後に移動するステップをさらに含み、前記第 2 のデータセンタは、前記優先リストにおいて前記第 1 のデータセンタよりも上位にあることを特徴とする請求項 8 に記載の方法。

【請求項 13】

前記複数のデータセンタは、ネットワークを通じてユーザがアクセス可能なネットワー

10

20

30

40

50

クサービスシステムを実現し、前記ネットワークサービスシステムは、ユーザがアクセス可能なウェブサイトを提供することを特徴とする請求項 8 に記載の方法。

【請求項 14】

前記ウェブサイトは、商業ウェブサイトであり、前記データセットは、ユーザのうちの一人のためのショッピングカートについてのショッピングカートデータを含むことを特徴とする請求項 13 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

この出願は、2006年3月31日に出願され、参照することによってここに組み込まれる「高可用性データを提供するためのシステム及び方法」という名称の米国特許出願第 11/394,648 号の一部継続出願である。

【背景技術】

【0002】

事業コンピュータ環境は、特定のビジネスアプリケーションに関連するデータにアクセスすることをしばしば必要とする。単一障害点 (single point of failure) を回避するために、データは、しばしば、異なるロケーション (例えば、所定のデータセンタ、異なるデータセンタ等における異なるロケーション) の複数のホストに記憶される。従って、例えば、特定のデータセットが一つのホストから利用不可能になったとしても (例えば、ホスト障害のために、ネットワークパーティション又は他のネットワーク障害等のために)、クライアントプロセスは、他のホストの当該データにアクセスすることができる。個々のホストは高可用性ではないこともあり得るが、個々のホストの組合せは、より高度の可用性の解決策を提供する。

【0003】

同一データを複数のロケーションに記憶する場合に直面する問題は、当該データの種々のコピー間における一貫性の保持である。一つのホストに存在するままのデータセットの状態は、他のホストに存在するままのデータセットの状態と一致しないことがある。例えば、クライアントプロセスが一つのホストにおけるあるデータセットに変更を発生させたとして、当該データセットはその後当該ホストからは利用不可能となり、当該ホストにおける当該データセットのコピーになされた変更は、少なくとも一時的に喪失されることがある。当該データセットの最近のバージョンは、別のホストから得られる可能性がある。しかしながら、クライアントプロセスが他のホストからの当該データセットにおける動作を開始すると、他のデータセットに反映されていない変更をそれぞれ伴う当該データセットの二つのバージョンが潜在的に生成され得るという更なる問題が発生する。

【0004】

従って、高可用性データを提供することが可能なシステム及び方法に対する継続的な要求が存在する。ある特徴及び利点が議論されているとはいえ、ここでの教示は、それらの特徴及び利点のいずれも必ずしも達成しないシステム及び方法を達成するためにも適用され得るということは理解されるべきである。

【図面の簡単な説明】

【0005】

【図 1】実施の一形態に係るデータ処理システムのブロック図である。

【図 2】実施の一形態に係る他のデータ処理システムのブロック図である。

【図 3】実施の一形態に係る図 1 のデータセットサービスを詳細に示すブロック図である。

【図 4】実施の一形態に係る図 1 のシステムにより実行される書込動作のフローチャートである。

【図 5】実施の一形態に係る図 1 のシステムにより実行される読出動作のフローチャートである。

10

20

30

40

50

【図 6】実施の一形態に係る図 1 のシステムにより実行されるデータ調整及び更新動作のフローチャートである。

【図 7】実施の一形態に係る図 1 のシステムにおけるデータ複製及び負荷平衡との関連において用いられるハッシュ動作の図である。

【図 8】実施の一形態に係る図 7 に示されるハッシュ動作の他の図である。

【図 9】実施の一形態に係る図 1 のシステムの増加スケーラビリティ機構を示す図である。

【図 10】実施の一形態に係る図 1 のシステムにおいて用いられるデータ複製構成の図である。

【図 11】実施の一形態に係る図 1 のシステムにおいて用いられるホスト優先リストの図である。 10

【図 12】実施の一形態に係る図 1 のシステムにおいて用いられる負荷平衡構成の図である。

【図 13A】実施の一形態に係る図 1 のシステムにより実行される書込動作のフローチャートである。

【図 13B】実施の一形態に係る図 1 のシステムにより実行される書込動作のフローチャートである。

【図 14A】実施の一形態に係る図 1 のシステムにより実行されるハンドオフ (hand-off) 動作を含む書込動作のフローチャートである。

【図 14B】実施の一形態に係る図 1 のシステムにより実行されるハンドオフ (hand-off) 動作を含む書込動作のフローチャートである。 20

【図 15A】実施の一形態に係る図 1 のシステムにより実行される読出動作のフローチャートである。

【図 15B】実施の一形態に係る図 1 のシステムにより実行される読出動作のフローチャートである。

【図 16】実施の一形態に係る図 1 のシステムにおいて用いられるデータバージョンング (versioning) 構成の図である。

【図 17】実施の一形態に係る図 1 のデータセットサービスを詳細に示すブロック図である。

【図 18】実施の一形態に係る図 17 のシステムにおけるデータ複製及び負荷平衡との関連において用いられるハッシュ動作の図である。 30

【図 19】実施の一形態に係る図 17 のシステムにおいて用いられるデータセンタ及びホスト優先リストの図である。

【図 20】実施の一形態に係る図 17 のシステムにより実行されるアクセス動作のフローチャートである。

【図 21】実施の一形態に係る図 20 のアクセス動作の一局面を詳細に示す図である。

【図 22】実施の一形態に係る図 20 のアクセス動作の一局面を詳細に示す図である。

【図 23】実施の一形態に係る図 20 のアクセス動作の一局面を詳細に示す図である。

【図 24】実施の一形態に係る図 20 のアクセス動作の一局面を詳細に示す図である。

【図 25】実施の一形態に係る図 17 のシステムにおいて用いられるメッセージフィルタ 40 である。

【発明を実施するための形態】

【0006】

実施の一形態は、複数のデータセンタ内のホストにデータセットを記憶することを含むコンピュータ実施 (computer-implemented) データ処理方法に関する。データセンタ及びデータセンタ内のホストは、例えば、複数段 (multi-tiered) リング構成に従って構成され得る。実施の一形態において、ハッシュ構成は、データセットの書込及び読出が発生するデータセンタ及びホストを選択するためのリング構成を実施するために用いられる。他の実施の形態において、バージョン履歴は、ホストにおける書込及び読出も行われ、また、バージョン履歴は、読出が発生した後のデータセット間の一時的な関係の評価に用いら 50

れる。

【0007】

発明の詳細な説明及び特定の実施例は、本発明の好適な実施の形態を示すと同時に、限定ではなく例示として与えられるということは理解されるべきである。本発明の範囲内において、その精神から離れることなく、多くの変形及び変更が可能であり、本発明は、総てのそのような変形を包含する。

【0008】

I. システム構成 (System Architecture)

図1を参照すると、実施の一形態に係るデータ処理システム100が示されている。データ処理システム100は、ユーザコンピュータ102と、通信ネットワーク104と、ネットワークサービスシステム106とを含んでいる。ユーザコンピュータ102は、通信ネットワーク104を介してネットワークサービスシステム106にアクセスし得る。ネットワークサービスシステム106は、ネットワークインタフェース110と、データセットサービス112と、一つ又は複数の他のサービス114とを含んでいる。ネットワークインタフェース110は、通信ネットワーク104を介してユーザからデータを受信し、ユーザにデータを供給する。例えば、ネットワークインタフェース110は、データセットサービス112により保持されているデータセットへのアクセスに加えて、他のサービス114により生成され及び/又は保持されている他のデータへのアクセスを、ユーザコンピュータ102に提供し得る。

【0009】

データセットサービスは、データセットを記憶し得るデータ記憶システム118を含んでいる。データ状態は、システム106内におけるユーザインタラクション(ユーザ相互作用)に基づいて及び/又は他の変化に基づいて時間と共に変化し得る。ここで、用語「データセット」は、時間と共に変化し得る任意のデータに関連する。例えば、各データセットは、当該データセットから付加、除去及び/又は変更され得る一つ又は複数の項目を含み得る。データ記憶システム118は、システム障害(例えば、ホスト障害、ネットワーク障害等)の場合には、後述するように、データセットが高度の一貫性を有して利用可能(可用)に存続するように、高可用性の方法において情報を記憶するように構成されている。実施の一形態において、データ記憶システム118は、バークリー(Berkeley)データベーストランザクションデータ記憶システムを用いて実現されている。

【0010】

ここで図2も参照すると、図2は、データ処理システム100の他の実施例を提供している。図2の例においては、ネットワークサービスシステム106は商業ウェブサイトシステム116であり、ネットワークインタフェース110はネットワークショッピングインタフェース120である。商業ウェブサイトシステム116は、例えば、数千又はそれ以上のホストを含む分散型計算方式環境において実現され得る。商業ウェブサイトシステム116は、アイテム(例えば、商品、サービス、予約申込等)の買い物をするためにユーザコンピュータ102を操作するユーザにとってアクセス可能である商業ウェブサイト(例えば、オンライン小売ウェブサイト)を提供し得る。そのような実施の形態においては、ネットワークショッピングインタフェース120は、アイテムの表示及び/又は販売を容易にするために、ウェブサイト上のグラフィックデータ及び/又はテキストデータをユーザに提供する。ユーザに提供されるデータは、価格、寸法、利用可能性、購入のために現在選択されているアイテム等のアイテム情報を含み得る。商業ショッピングインタフェース120は、ユーザが関心を持っているアイテムを示すデータ、トランザクションの完了を必要とするデータ等、ユーザからのデータを受信するように構成されることも可能である。

【0011】

図2の実施例においては、データセットサービス112は、ウェブサイトのユーザによって購入のために又は購入可能性のために選択されたアイテムのリストを保持するショッピングカートデータサービス122となるべきものとして示されている。そのような実施

例においては、各データセットは、特定の顧客に関連するショッピングカートであり得る。データセットは、ショッピングカート内のアイテムのためのアイテム識別情報、ユーザが選択したかもしれないがまだ購入していないアイテムのためのアイテム情報、ショッピングカート内のアイテムの数量情報等を含み得る。ショッピングカートデータサービス122は、ショッピングカートに関連する他のビジネスロジックを含み得るショッピングカートサービス124を通じてアクセスされ得る。ウェブサイトシステム116は、例えばユーザのショッピングカートの全部又は一部を表示するウェブページ等、データセットの全部又は一部を含むウェブページを、ウェブサイトのユーザのために表示し得る。他の例示的な実施の形態においては、データセットは、ユーザのインタラクション（相互作用）に基づいて、又は、訪問者の便宜のために、又は、ウェブサイトの操作を容易にするために、ウェブサイトシステム116により収集され得る他のデータを含み得る。例えば、データセットサービス112は、特定のエンティティ（構成要素、実体）に関連するデータセット（例えば、ウェブサイトの異なるユーザ、ウェブサイト上の異なるセッション、ウェブサイト上で実施される異なるトランザクション、ウェブサイトによって提供される異なるアイテム、ウェブサイトによって提供されるアイテムの異なるカテゴリ、ウェブサイト上で表示される異なる広告、ウェブサイトの異なるページ等に関連するデータセット）も保持し得る。図2はウェブサイトシステムを示しているが、理解されるであろうように、データ処理システム100は、他のアプリケーションにおいても使用され得る。

【0012】

再度図1を参照すると、データセットサービス112は、ローカル（局所）プロセス及びリモート（遠隔）プロセスの両方との関連において使用され得る。リモートプロセスとの関連においては、データセットサービス112に対する読出及び書込要求は、通信ネットワーク104経由でリモートプロセスから受信され得る。例えば、ネットワークサービスシステム106は、アプリケーションプログラムインタフェース（API）を介してインターネットを通じてリモートプロセスにアクセス可能なサービスを提供し得る。そのようなサービス要求は、サードパーティによって、例えばそれら自体のデータ処理システムの動作における援助（アシスト）を行うために、行われ得る。

【0013】

ここで図3乃至図6を参照すると、データセットサービス112の構成及び動作が詳細に示されている。図3に示されるように、データセットサービス112は、複数のホスト130を含み得る。ここで、用語「複数」は、二つ又はそれより多いことを意味する。例えば、データセットサービス112は、数十、数百若しくは数千又はそれより多いホストを含み得る。実施の一形態において、各ホスト130は、機能的に等価である（例えば、同一コードを実行し、又は、同一コードの関連するバージョンを実行する）。各ホスト130は、後述する図3乃至図16に記載された動作を実行するように構成された、格納されたプログラムロジックを含み得る。以下に述べられるように、データセット記憶システム118は、各ホスト130がデータセットの一部を記憶するように、各ホスト130に亘って分散される。各ホスト130は、（キー値ペア（key-value pairs）の）データのサブセットを記憶し、システムは、各データセットのNの複製（ここで、Nは、複製係数、又は、データセットを複製する回数を表す正の整数である）を保持することを試みる。値Nは、設定可能であり、データの耐久性、可用性及び一貫性のいずれにも作用する。システム内にSの物理的なホストがあるとすると、全システム106は、S・Nの物理的なホストを含み（Sが小さいほど全システムの可用性は低くなるが）、各ホスト130は、データセットの約N/Sを記憶する。代替的に、異種のホスト130が使用される場合、各ホスト130は、システム106内に重み付けをする各ホスト130の重みに比例するいくつかのデータセットを記憶する。各ホスト130の重みは、各ホスト130の資源に基づいて決定され得る。例えば、各ホスト130の重みは、より高性能なホスト130がより多くのデータセットを記憶するように、各ホスト130の相対的な性能に基づいて（例えば、処理能力、記憶容量、及び/又は、ネットワーク容量に基づいて決定されるように）決定され得る。Nの値は、例えば、データセットごと又はデータタイプごとの基準で

10

20

30

40

50

可用性 / 耐久性が設定されることを許容するために、データセットごと又はデータタイプごとの基準で設定可能なものとされ得る。

【 0 0 1 4 】

図 4 に示されるように、クライアントプロセス 1 3 4 (例えば、サービス 1 1 4 のうちの一つ) から受信されるデータを記憶するために、データセットサービス 1 1 2 は、クライアントプロセス 1 3 4 から書込要求を受信し (ステップ 1 5 0)、その後、複数のホスト 1 3 0 にデータ書込を行うことによって応答する (ステップ 1 5 2)。(この出願の目的のため、用語「クライアントプロセス」とは、任意の他のプログラムロジックからの、例えばここではデータセットサービス 1 1 2 からの、データセットを要求することがある任意のプログラムロジックをいう。)実施の一形態においては、後述するように、データは、優先リストに基づいて複数のホストに書き込まれる。データが書き込まれた後、書込動作が実行されたことを確認する応答が、クライアントプロセス 1 3 4 に送信される (ステップ 1 5 4)。例示的な書込動作が、図 7 乃至図 1 2、図 1 3 A 乃至図 1 3 B 及び図 1 4 A 乃至図 1 4 B との関連において詳細に記載されている。

【 0 0 1 5 】

図 5 に示されるように、クライアントプロセス 1 3 4 にデータを供給するために、データセットサービス 1 1 2 は、クライアントプロセス 1 3 4 から読出要求を受信し (ステップ 1 6 0)、その後、複数のホスト 1 3 0 においてデータ読出を行うことによって応答する (ステップ 1 6 2)。データが読み出された後、読出動作が実行されたことを確認し且つ要求されたデータを含む応答が、クライアントプロセス 1 3 4 に送信される (ステップ 1 6 4)。例示的な読出動作が、図 1 5 A 乃至図 1 5 B との関連において詳細に記載されている。

【 0 0 1 6 】

図 6 に関しては、総ての関連するネットワーク接続及びホスト 1 3 0 が健全である (例えば、利用可能であり応答する) 場合、読出動作に係るホスト 1 3 0 は、典型的には、一貫性のあるデータを供給する。しかし、一つ又は複数のネットワーク接続又はホスト 1 3 0 に故障又は障害が発生している場合、ホスト 1 3 0 は、同一データセットの異なるバージョンを供給し得る。従って、図 6 に示されるように、クライアントプロセスにおいてデータセットが受信された後 (ステップ 1 7 0)、データセットは、調整されて一致させられ得る (ステップ 1 7 2)。調整されて一致させられたデータセットは、その後、記憶のためにデータセットサービス 1 1 2 に送信され得る (ステップ 1 7 4)。以下に詳細に記載されるように、同一データセットの一貫性のないバージョンの存在は、データバージョンニング (data versioning) 構成を用いて削除され得る。データバージョンニング構成は、一貫性のないバージョンを調整するために、バージョン調整ロジック 1 3 6 (図 3 に示されるように、クライアントプロセス 1 3 4 の一部として又はクライアントプロセス 1 3 4 との関連において提供される) によっても使用され得る。例示的なデータバージョンニング構成が、図 1 6 との関連において以下に詳細に記載される。

【 0 0 1 7 】

I I . 読出 / 書込動作の協調 (Coordination of Read/Write Operations)

A . ホスト間のデータセットの区分化 (Partitioning Data Sets over Hosts)

図 7 乃至図 8 を参照すると、実施の一形態において、データセットサービス 1 1 2 は、システム 1 0 6 内のホスト間にデータセットを区分化するための機構を含んでいる。実施の一形態において、以下に記載されるように、一貫性ハッシュ構成は、データがホスト 1 3 0 間に相対的に均等に分散されるように、データセットを記憶するために使用され得る。他の実施の形態においては、他のデータ区分化構成が使用され得る。

【 0 0 1 8 】

先ず図 7 を参照すると、実施の一形態において、データセットサービス 1 1 2 によって記憶されているデータにアクセスするために (例えば、読出動作又は書込動作を介して)、クライアントプロセスは、各要求が参照する、データセットのためのキーを含むデータ要求を送信する。例えば、ショッピングカートアプリケーションとの関連においては、キ

ーは、ショッピングカートが関係するユーザのユーザIDに基づいて生成され得る（例えば、ユーザIDがキーとして使用され得る）。キーは、データセットに関連し且つハッシュ関数への入力としての使用に適当な任意のデータ値であり得る。図7に示されるように、キーは、キーの関数としてのハッシュ値hを順次生成するハッシュ関数182に適用される。実施の一形態において、ハッシュ関数182は、ハッシュ範囲に亘るハッシュ値のほぼ均一な分散を達成する。図示された実施の形態においては、ハッシュ値は、ハッシュ範囲 $\{0, 2^{128}\}$ に亘って分散するように示されているが、ハッシュ値の任意の数、又は、実際上は任意の大きさのハッシュ範囲が使用され得る。

【0019】

データセットサービス112においてアクティブな（活動状態にある）参加者になると、各ホスト130は、ハッシュ範囲に亘る位置のセットを割り当てられる。説明の目的のために、ここでの考察の残りの部分では、データセットサービス112を実施するホストが、ホストA、ホストB、ホストC、ホストD及びホストEとして示されている五つのホスト130であると仮定する。実際には、データセットサービス112が、数十、数百若しくは数千又はそれより多いホスト130により実施され得ることは、理解されるであろう。

【0020】

図8を参照すると、図8は、読出動作又は書込動作に対する責任（分担）がハッシュ値に基づいて特定のホスト130に割り当てられる方法を示している。各ホスト130は、ハッシュ範囲内におけるそれ自体の位置から先行するホスト130の位置までに亘るハッシュ値との関連において読出／書込動作に対して責任を有する。例えば、ホストA、B、C、D、Eがハッシュ値 h_1, h_2, h_3, h_4, h_5 にそれぞれ位置しているとする、ホストBは、ハッシュ値 h_1, h_2 の範囲に対して責任を有し、ホストCは、ハッシュ値 h_2, h_3 の範囲に対して責任を有し、以下、同様となる。ホストAに対する責任（分担）「ラップアラウンド（wrap around）」の割当は、即ち、ホストAが、ハッシュ値 $h_5 < h_2^{128}$ 及び $0, h_1$ の範囲に対して責任を有することである。実施において、例えば、キー k_1 及び k_2 を有するデータセットは、キー k_1 及び k_2 をハッシュすることによりリング184上におけるそれらの位置を与えて、ホスト130に割り当てられ、その後、リング184を時計回りに周回して、データセットのハッシュされたキーよりも大きい値の位置を有する最初のホスト130を見出す。キー k_1 の場合、対応するデータセットが割り当てられる、より大きい位置を有する最初のホストは、ホストAである。キー k_2 の場合、対応するデータセットが割り当てられる、より大きい位置を有する最初のホストは、ホストBである。

【0021】

図7乃至図8に示される構成は、各ホスト130がリング184上におけるそれ自体とその先行ホスト130との間のリング184の範囲に対して責任を有するという結果になる。例えば、ホストBは、それ自体とホストAとの間のリング184の部分に対して責任を有する。ある一つのホスト130が参加又は離脱したとすると、そのことは、リング184上におけるその直近の後続ホストの責任（分担）に影響を与えるだけであり、他の総てのホスト130は影響を受けない。これが図9に示されており、そこでは、ホストFの追加が、リング184上におけるその直近の後続ホストであるホストBの責任（分担）に影響を与えるが、ホストA等の他のホストの責任（分担）には影響を与えない。従って、個々のホスト130は、各ホスト130に対するデータセットの区分化の全体的な再配置を伴うことなく、追加又は除去されることが可能であり、それにより、増加スケーラビリティ（拡大縮小可能性）（scalability）が増進される。

【0022】

1. データ複製（Data Replication）

ここで図10乃至図11を参照すると、図7乃至図8のハッシュ構成は、データ複製の支援のために使用され得る。図10においては、リング184上における直近の後続ホスト130に対して単にデータセットが割り当てられるよりむしろ、データセットは、最初

10

20

30

40

50

のNの後続ホスト130に割り当てられている。後述するように、データセットサービス112は、ホスト130間にデータのNの複製があり、各ホスト130がそれ自体とそのN番目の先行ホストとの間のリング184の範囲に対して責任を有するということを保証するために機能する。

【0023】

図11に示されるように、このような構成においては、各キーは、当該キーに基づいて生成されたハッシュ値からリング184を時計回りに周回したときにホスト130のそれぞれが最初に遭遇を受ける順序である、ホスト130の優先リスト190を有している。優先リスト190は、データセットにアクセスするために（例えば、読出又は書込のために）使用されるホスト130の優先順序を表している。総てのホスト130が健全である場合、優先リスト190における上位Nのホストがデータセットを記憶する。特定のホスト130に障害が発生した場合、又は、ネットワークパーティションが発生した場合、データセットは、優先リスト190においてより低いランク付けのホスト130に一時的に記憶され得る。複数のホスト130に障害が発生した場合、データセットは、優先リスト190においてより低いランク付けの複数のホスト130に記憶され得る。N=3のとき、キー k_1 に関連するデータセットにアクセスしているクライアントプロセス134は、キー k_1 の位置からリング184を時計回りに周回することによって分かるように、ホストA、B及びDに（より上位のホストのいずれかが利用不可能である場合には次にホストEにさらにホストCに）この順序で読出又は書込を行う。キー k_2 に関連するデータセットにアクセスしているクライアントプロセス134は、キー k_2 のハッシュ位置からリング184を時計回りに周回することによって分かるように、ホストB、D、Eに（より上位のホストのいずれかが利用不可能である場合には次にホストCにさらにホストAに）この順序で読出又は書込を行う。以上に述べたように、値Nは設定可能な値であり、従って、データセットのより多くの複製を許容するために、より多くのホスト130がシステム106に付加され得る。従って、データセットの可用性のレベルは設定可能であり、適当な数のホスト130を使用して所望するだけの高さに設定され得る。

【0024】

総てのホスト130が利用可能である場合、同一データセットにおける連続する動作はNのホストの同一セットにアクセスし、従って一貫性を有する（即ち、動作は、同一キーにおける先行動作により読出/書込が行われた同一データにアクセスする）。ネットワーク又はホストの障害がある場合、同一データセットへの連続する動作は、ホスト130の異なるセットにアクセスすることがあるが、動作は、アクセスされるホストのセットにいくらかの重複がある限り、依然として一貫性を有し得る。例えば、キー k_1 における第1の動作が、ホストA、B及びDにアクセスし得る。その後、ホストBが利用不可能になった場合、キー k_1 における第2の動作は、ホストA、D及びEにアクセスし得る。従って、優先リスト190における最上位の利用可能なホスト130にアクセスすることにより、動作から動作へのホストの利用可能性の小さい変化は、一貫性に否定的な影響を与えない。その理由は、後続のアクセスが、重複するホストを含み得るからである。（非一貫性という結果になる）ホストセットの間に重複が存在しない状態となるためには、少なくともNのホストの利用可能性（可用性）は、二つの連続する動作の間に変化しなければならない。以上に述べたように、値Nは設定可能な値であり、従って、一貫性の確率保証は設定可能であり、所望するだけ高く設定し得る。これは、大域（global）一貫性（システム応答は、データに対してなされた最新の絶対的な変更を反映する）、及び、主観的一貫性（システム応答は、現在の要求を行っているクライアントによりなされた最新の変更を反映する）の両方の確率保証を含む。

【0025】

実施の一形態において、データセット上におけるクライアント動作は、複数のロケーション（例えばサーバ）において提供され得る。また、同一データセット上における連続する動作は、異なるサーバによって提供され得る。実施の一形態において、所定のデータセットを記憶するホスト130にアクセスするために、サーバは、ホスト130の利用可能

10

20

30

40

50

性（優先リスト190において最上位であるNの利用可能なホストを選択するために）に加えて、ハッシュ空間におけるホスト位置に関する情報（優先リスト190を計算するために）を記憶する。ネットワーク又はホストの障害が存在する際には、異なるサーバが、ホストの利用可能性に関する異なる情報を記憶し得る。システムに結合され又は離脱するホストが存在する際には、異なるサーバが、ハッシュ空間におけるセット位置に関する異なる情報を記憶し得る。例えば、サーバXは、ホストAがデータセットサービス112に結合したことを認識しないことがある。従って、キー k_1 を有するデータセット上における動作の提供において、サーバXは、ホストB、D及びEにアクセスし得る。別のサーバYは、ホストAとホストAのハッシュ位置とを既に認識している可能性がある。この情報に基づいて、キー k_1 上における後続の動作を提供する際には、サーバYは、ホストA、B及びDにアクセスし得る。従って、優先リスト190における最上位の利用可能なホスト130にアクセスすることにより、書込及び読出動作の間に少なくとも一つのホストにアクセスする確率は増加する。以上に述べたように、一貫性のこの確率保証は、Nの値によって決定される。

【0026】

実施の一形態においては、優先リスト190は、ハッシュ関数182の動作によって実施され得る（例えば、分離して記憶されることなく）。他の実施の形態においては、優先リスト190は、記憶され得る。理解されるであろうように、他の要因が、優先リスト190を構成する際に考慮に入れられることがある。優先リスト190は、そのような要因を考慮に入れるために、手動で又は自動的に構成され得る。例えば、可用性及び耐久性をさらに改善するために、同一優先リスト190内に相互関係障害の確率が比較的低いホスト130が含まれるように、優先リスト190が構成され得る。例えば、システム100が複数のネットワークに亘って分散させられている場合、同時に故障することがありそうもないホスト130の集合があり得る。従って、システム100は、データセットのNの複製のために、障害の相互関係が低くなるようにNのホストを選択することによって、可用性及び耐久性を最大化することができる。同様に、低い障害相互関係は、ホスト130が異なるハードウェア上で運転され、異なるプログラムロジック実行を用い、地理的に異なる領域において運転され、それらが組み合わせられる場合にも存在し得る。例えば、リング184を時計回りに周回するときに、遭遇したホスト130が、考慮されるべきことが望ましい任意の追加的な標準（criteria）に適合するか否かの評価に、規則の集合が適用され得る。遭遇したホスト130が追加的な標準に適合しない場合、利用可能なホストの検索が、追加的な標準に適合するホストに遭遇するまで、リング184の周回を前進して継続され得る。

【0027】

他の構成も、地理的な多様性を実現するために用いられ得る。そのような構成の実施例が、図17乃至図25との関連において以下に詳細に記述される。

【0028】

2. 負荷平衡 (Load Balancing)

図12を参照すると、ホスト130は、負荷平衡を促進するために、即ち、データ及び負荷の不均一な分散を回避するために、リング184上の複数の位置に割り当てられ得る。そうでなければ、リング184上における各ホスト130の無作為な位置割当によって、データ及び負荷の不均一な分散が発生し得る。従って、図12において、ホストA、B、C、D、Eは、リング184上の複数の位置に割り当てられている。この複数の位置決めは、各ホスト130に割り当てられたデータセットの数の分散を減少させる傾向がある。その理由は、リング184上における増加した多数の無作為の配置が、各ホスト130に割り当てられたデータセットの数の平均値への集中を引き起こす傾向があるからである。従って、より多くの位置をリング184上における各ホスト130に割り当てることは、負荷平衡を改善する。実施の一形態において、遭遇した各ホスト130の最初の実例だけが優先リスト190に配置される。キー k_1 の場合、対応するデータセットが割り当てられている、より大きい位置の最初のホストは、ホストAである。N = 4の場合、キーk

k_1 に関連するデータセットにアクセスするプロセスは、ホストA, B, C, Dに対して読出又は書込を行う。キー k_1 のための優先リスト190は、リング184上における複数の位置を有するホストのために、及び、異なる順序で遭遇するホストのために、上記とは異なる。キー k_2 の場合、対応するデータセットが割り当てられている、より大きい位置の最初のホストは、ホストBである。キー k_2 に関連するデータセットにアクセスするクライアントプロセス134は、ホストB, C, D, Aに対してその順序で読出又は書込を行う。他の例示的な実施の形態においては、遭遇する各ホスト130の複数の実例は、例えば、以前は利用不可能であったホスト130を再試行するために、優先リスト190に配置され得る。

【0029】

10

ホスト130をリング184上の複数の位置に割り当てることは、異種のハードウェアの使用を容易にし、即ち、より性能の高いホスト130がリング184上のより多くの位置に割り当てられ、より性能の低いホスト130がリング184上のより少ない位置に割り当てられ得る。例えば、図12において、ホストEは、他のどのホストよりも少ない位置を有しており、従って、性能の低いホストであると想定される。理解されるであろうように、あるホストの範囲が使用されることがあり、各ホストは、他のホスト130より性能が高いか又は低い。特定のホスト130に割り当てられている位置の数は、当該特定のホスト130の相対的な性能の関数であり得る。

【0030】

加えて、十分な数の位置が各ホスト130に割り当てられたとすると、各ホスト130は、他のホスト130のそれぞれと後続/先行関係を有し得る。従って、ホスト130のうちの一つが利用不可能になり又は使用解除されたとすると、使用解除されたホストにより処理されていた負荷は、データ可用性を失うことなく残余の利用可能なホスト130に亘ってほぼ均等に分散させられ得る。同様に、あるホスト130が再度利用可能になった場合又は新たなホスト130がデータセットサービス112に追加された場合は、その新たに利用可能になったホスト130は、他の利用可能なホスト130のそれぞれからおおよそ均等な量の負荷を軽減し得る。

20

【0031】

B. 読出/書込アクセス動作 (Read/Write Access Operations)

ここで図13A乃至図13B、図14A乃至図14B及び図15A乃至図15Bを参照すると、読出及び書込動作が示されている。読出/書込動作は、クライアントプロセス134によりデータセットサービス112に対して行われるサービス要求によって引き起こされ得る。サービス要求を受信すると、データセットサービス112は、要求された動作を実行し、クライアントプロセス134に対して応答を供給する。

30

【0032】

データセットサービス112においては、ホスト130のうちの一つが、読出又は書込要求を調整するための責任を有する。読出又は書込要求を調整するための責任を有するホスト130を、ここではコーディネータ(調整者)と称する。実施の一形態において、コーディネータは、優先リスト190に載せられた最初のホスト130であり、ローカル読出又は書込動作の実行を含む読出又は書込要求を調整する。例えば、サービス要求は、最初は別のホスト130によって受信されることがあり、当該ホスト130は、コーディネータとして機能するホスト130(例えば、優先リスト190における最上位ホスト)にそのサービス要求を転送することを決定し得る。他の実施の形態においては、コーディネータは、優先リスト190に載っていないホスト130等の他のホスト130であり得る。例えば、コーディネータは、偶然に最初に読出又は書込要求を受信したが、たまたま優先リスト190の最上位付近ではなく、優先リスト190の最上位付近にあるホストにそのサービス要求を転送することを決定しないホスト130であり得る。実施例を提供する目的のために、ここでは、コーディネータは優先リスト190に載せられた最初のホスト130であるものとする。

40

【0033】

50

実施の一形態においては、上述のように、読出及び書込動作は、潜在的にダウンしている又はアクセス不可能なホスト130をスキップしながら、優先リスト190における上位Nの健全なホストにアクセスし得る。総てのホスト130が健全である場合は、あるキーの優先リスト190における上位Nのホストがアクセスされ得る。ホスト障害又はネットワークパーティションが存在する場合、優先リスト190におけるより下位のホスト130が代わりにアクセスされることがあり、それにより高可用性が維持される。

【0034】

先ず図13A乃至図13Bを参照すると、例示的な書込動作が示されている。図13Aにおいて、バージョン V_{n+1} のための書込要求が、ホストAによってクライアントプロセス134から受信される（上述のように、直接に又は間接に、のどちらでも）。図12に示されるようなリング184上におけるホスト130の分散を想定すると、キー k_1 のための優先リスト190は、 $P_L = \{A, B, C, D, E\}$ である。この実施例において、ホストAはコーディネータであり、書込動作をローカルに実行する（ステップ150）。ホストAは、その後、新しいバージョン V_{n+1} を残余の最上位にランクされているN-1の到達可能なホストであるホストB及びCにコピーし（例えば、 $N=3$ のとき）、それはその後、書込動作の実行及び追加的なコピーの記憶も行う（ステップ152）。

【0035】

データセットが記憶されるとき、データ自体に加えて、当該データに関連するキー、及び、ベクトルクロックも記憶される。キーは、データセットが後に識別されることを可能にする。ベクトルクロックは、同一データセットの異なるバージョン間の因果関係を捕捉するためのデータバージョニングのために使用され、データセットのバージョンに関連する{ホストID, カウンタ}ペアのリストを含む。ベクトルクロックの使用を通じたデータバージョニングは、図16との関係において詳細に後述される。

【0036】

図13Bにおいては、ホストB及びCは、書込動作が成功したか否かをホストAに折り返し報告し、ホストAは、書込動作が成功したか否かの確認をクライアントプロセス134に応答する（ステップ154）。実施の形態において、成功と考えられるべき書込動作のために、書込動作は、Wが設定可能な値であって $W \leq N$ であるとする、Wのホストにおいて成功しなければならない。従って、例えば、 $N=3$ かつ $W=2$ とすると、書込動作は、三つのホスト130において試行されたとしても、二つのホスト130において成功すれば、成功であると考えられる。書込動作が一つ又は複数のホスト130において成功したとすると、データセットのコピーは、詳細に後述するように、依然としていつかは優先リスト190における上位Nのホストに移動し得るということを意味し得る。従って、上述したテストに従って書込動作が成功であると考えられなかったとしても、上位Nのホストにおけるデータセットの最終的な一貫性は依然として実現され得る。

【0037】

図14A乃至図14Bを参照すると、データハンドオフを伴う例示的な書込動作が示されている。データハンドオフは、あるデータセットのための優先リスト190における最上位にランクされているNのホストに対してデータを移動させることを試行する機構又は手順である。例えば、上述のように、コーディネータは、一般に、優先リスト190における上位Nのホストに対してデータの送信を試行する。しかしながら、一つ又は複数のホスト130がダウンしている場合、コーディネータは、優先リスト190におけるより下位のホスト130に対してデータを送信する。優先リスト190は、書込動作に（及び読出動作に）関係することとなるホスト130の適切に規定された順序を提供し、データハンドオフ機構は、優先リスト190における最上位にランクされているNのホストに対してデータを移動させて返すために使用される。

【0038】

従って、図14Aに示されるように、ホストAは、図13Aのようにバージョン V_{n+1} のための書込要求を受信する。ホストAは、その後、書込動作を実行し、残余の最上位にランクされているNの到達可能なホストであるホストB及びCに新しいバージョンをコ

10

20

30

40

50

ピーすることを試行する。図示された実施例においては、ホストCが一時的に故障しており、従って、ホストDにおける書込が試行される。ホストDにおいて書き込まれたデータは、そのいくらか後の時点でホストDが当該データをホストCに転送し得るように、どのホスト130がデータを受信して書込を行うべきか（例えば、ホストC）を提案するヒントと共にタグを付され得る。図14Bにおいて、ホストCが健全である場合、データハンドオフが行われ、データはホストCにコピーし戻される（copied back）。データは、このようにして、優先リスト190における最上位にランクされたNのホストの一つであるホストCに移動し戻される。

【0039】

実施の一形態においては、関連する技術がデータセットの喪失したコピーを回復するために使用され得る。例えば、ホスト130が参加又は離脱して優先リスト190に対応する変更があるとすると、データの不適切な配置が発生し得る。例えば、システム100に追加されたホスト130は、優先リスト190における他のホストの序列に不適切な配置を発生させる。そのような状況において、データハンドオフを実行するために、ホスト130のペアは、それらが共通に共有する領域の比較を定期的に行い、その後、比較の間に検出されたいかなる差異も調整して一致させるために、必要なデータ転送を実行し得る。例えば、キーの領域を自体のために保持している、上位Nのホストの一つではないホスト（送信者）は、上位Nのホストの任意の一つ（受信者）を無作為に選択し得る。別の実施例として、上記ホストは、上位Nのホストのなかから、例えば、当該データを有している可能性の低いホストを選択し得る。その理由は、上記ホストがデータセットサービス112に最近参加したものである。上記二つのホスト130は、その後、低レベルデータベース比較を継続してもよく、送信者は、比較によって検出されたいかなる差異も調整して一致させるために、受信者が記憶しているものよりも新しいデータセットを転送し得る。データは、優先リスト190における少なくとも一つのホスト130に移動されることが可能であり、従って、優先リスト190における残余のホスト130に増殖される。例えば、残余のホスト130への増殖は、いくつかのキーのセットに対する優先リスト190における上位Nのホストに含まれているホスト130のペアに記憶されているデータセットを比較することによって実行され得る。実施の一形態においては、二つのホストに記憶されているデータ間の差異を効果的に見出すために、マークルツリー（Merkle tree）が使用され得る。例えば、マークルツリーは、ツリーの各ノードがそのサブツリー（subtree）におけるデータから算出された合計（又はハッシュ値）を含んでいるところで、及び、リーフ（leaves）が一つ又は複数のデータ値（例えば、キー、バージョン及びクロック）のハッシュを含んでいるところで、使用され得る。ツリーのコンテンツにおける差異は、ブランチ（branch）に沿ってデータ合計（ハッシュ値）が異なっている当該ブランチを再帰的に展開する（recursing down）ことによって見出され得る。比較の効果を改善するために、マークルツリーは、ブルームフィルタ（Bloom filter）を使用してコード化（符号化）され得る。

【0040】

上述した機構を用いて、データセットサービス112は、データセットの最も新しいバージョンのコピーをその優先リスト190における上位Nのホストに動的に移動するために、進行中の試行を行う。従って、データセットの最も新しいバージョンのコピーが最初はその優先リスト190におけるより下位のホスト130にコピーされることがあったとしても、当該コピーは最終的には優先リスト190における上位Nのホストに移動し戻され、上位Nのホストにおけるデータセットの最終的な一貫性に帰着する。

【0041】

図15A乃至図15Bを参照すると、優先リスト190を用いて実行される例示的な読出動作148が示されている。図15Aにおいて、読出要求は、ホストAによってクライアントプロセス134から受信される（上述のように、直接に又は間接に、のどちらでも）（ステップ160）。ホストAは、ローカル読出を行うためにホストB及びCからのデータを並行に要求することにより読出動作を調整する。ホストB及びCは、要求された読

出動作を実行する。図 15 B において、ホスト A は、ホスト B 及び C から読出結果を受信し（ステップ 162）、クライアントプロセス 134 に対して応答を供給する。

【0042】

読出要求を受信するとき、コーディネータは、そのキーに対する優先リスト 190 における最上位にランクされている N の到達可能なホスト 130 からの、そのキーに対するデータの総ての存在するバージョンを要求することがあり、その後、クライアントプロセス 134 に対して結果を戻す前に、R の応答を待つ（ここで、R は、良好な読出動作に参加することが必要とされるホストの数である）。図 15 A 乃至図 15 B の実施例においては、値 R は、3 に等しく設定される。

【0043】

値 W と同様に、値 R は設定可能である。例えば、 $R = 1$ であるとする、ホスト A は良好な読出について一度応答し、その読出からのデータは、クライアントプロセス 134 へ使用のために戻される。別の実施例として、 $R = 2$ であるとする、ホスト A 及び B の両方において読出が実行されるまでデータが戻されないことがある。二つの読出が実行されると、システム 100 は、データが同一バージョンであることを認識し、 $R = 1$ のときと同一のデータを戻す。さらに別の実施例として、 $R = 3$ であるとする、ホスト A、B 及び C において読出が実行されるまでデータが戻されないことがある。

【0044】

値 R 及び W は、一貫性及び高性能を提供するために、N より小さくなるように設定され得る。 $R + W > N$ となるような値 R 及び W の設定は、読出及び書込動作に関係するホスト 130 のセットの間における重複の設定可能な高い確率が存在するクォーラムのようなシステム（quorum-like system）をもたらす。より高い N が設定されると、少なくとも一つの複製が存在する可能性が高いため、システムは可用性及び耐久性を有する可能性が高くなる。一方、データは、ホスト 130 の同一セットに書き込まれ又はホスト 130 の同一セットから読み出される必要がないことは注目され得る。例えば、データセットは、優先リスト 190 におけるより下位のホスト 130 に書き込まれ、優先リスト 190 におけるより上位のホスト 130 にデータハンドオフを通じて移動されることがあり、その後、最終的には、優先リスト 190 におけるより上位のホスト 130 から読み出される。優先リスト 190 における上位 N のホストにおけるデータセットの可能性のある一貫性は、達成される。他の実施の形態においては、R 及び W は、N よりも十分に小さくなるように（例えば、 $R + W < N$ ）設定されることがあり、データセットのコピーは、（コーディネータに加えて） $W - 1$ のホストにのみ送信され得る。そのような実施の形態においては、上述のデータ回復機構は、上位 N のホストの残余のものにデータセットを増殖させるために使用され得る。

【0045】

実施の一形態において、データセットサービス 112 のためのアプリケーションプログラムインタフェースは、以下のように設定され得る。例えば、コマンドは、以下の形態を有し得る。

【0046】

`write(Key, Value, Context) ResultCode`

（書込（キー、値、コンテキスト） 結果コード）

`read(Key) ValueList, Context, ResultCode`

（読出（キー） 値リスト、コンテキスト、結果コード）

ここで、キー（Key）は、バイトの非有界シーケンスであり、値（Value）は、データ（バイトの非有界シーケンス）及びメタデータ（値が書き込まれた最終時刻を含む、値についての情報、ダイアグノスティック（診断）及びデバッグ情報等を包含する、読出専用の、任意の、拡張可能なデータセット）を含むオブジェクトであり、値リスト（ValueList）は、値（Value）のリストであり、コンテキスト（Context）は、読出 - 変更 - 書込サイクルに対するベクトルクロック状態をトラック（追跡）するために記憶システムによって内部構造において使用されるオpaque（opaque）オブジェクト（不透明なオブジェクト）で

10

20

30

40

50

あり、結果コード (ResultCode) は、読出又は書込動作が成功したか否かについてのコード表示である。

【 0 0 4 7 】

書込動作は、介在する書込がそのキーについて既に起こったことを意味するコンテキストの無効が発生していない限り、キーによって識別される値を、値パラメータによって特定される値に変更する。実施の一形態において、クライアントプロセス 1 3 4 は、読出 - 変更 - 書込サイクル (楽観的ロック (optimistic locking)) を再始動する。他の実施の形態においては、クライアントプロセス 1 3 4 は、データセットのコンフリクト (矛盾) するバージョンが存在し得る場合に、書込動作が継続することを許可し得る。読出動作は、キーに関連する値のための、データセットサービス 1 1 2 内におけるルックアップを実行する。良好に読み出された任意の及び総ての値は、値リストに戻される。オペークコンテキストオブジェクトは、後続の更新動作において使用するために戻される。複数の値が戻される場合、クライアントプロセス 1 3 4 は、総ての値に対して調整動作を実行することが期待される。後続の更新が (戻されたコンテキストを使用して) 実行される場合、更新された値が値リストに戻された総ての値の調整を表すという仮定は、(もしあれば) その値に任意の追加的な変更を加える。

【 0 0 4 8 】

理解されるであろうように、アプリケーションプログラムインタフェースにおける複雑性のより大きい又は小さいレベルが用いられ得る。例えば、実施の一形態において、値オブジェクトは、どの程度の期間だけデータが維持されるべきか、例えば、その結果、古い / 放棄されたデータが最終的に削除されてもよい、に関する情報が特定されることを可能にするタイプパラメータを含み得る。

【 0 0 4 9 】

他の実施の形態において、キーは、二つの部分 (パーティションキー、オブジェクトキー) に分割されて使用され得る。そのような実施の形態においては、パーティションキーは、キーパラメータに関して上述したように、そのキーに対する優先リスト 1 9 0 を生成するためにハッシュされ得る。同一のパーティションキーを共有する二つのデータセットは、従って同一の優先リスト 1 9 0 を有し、それ故に非常に高い確率でそれらのデータセットのそれぞれのコピーは、ホスト 1 3 0 の同一セットに属することとなる。パーティションキーを共有する総てのキーに対する優先リスト 1 9 0 の上位 N のホストのなかにホストの同一セットがあるので、そのような機構は、数個のデータセットに同時にアクセスすることを最適化として許容する。例えば、図 2 の商業ウェブサイトの例において、特定のユーザに関連する総てのデータセット (例えば、ショッピングカート、プロフィール、クレジットカード情報等) をホスト 1 3 0 の同一セットに記憶することは望ましいことであり得る。それらのデータセットのそれぞれについて同一のパーティションキーを使用することにより、データセットは、ホスト 1 3 0 の同一セットに記憶される。(パーティションキー、オブジェクトキー) 組合せは、ユーザに対する各個人のデータセットを一意に識別する。この構成によって可能となる他の最適化は、パーティションキーを共有するキーにおける範囲問合せ (range query) である。例えば、そのような範囲問合せは、所定のパーティションキーに対する優先リスト 1 9 0 の上位 N のホストのうちの単一ホスト 1 3 0 にアクセスすることによって、そのパーティションキーに対する総てのオブジェクトキーを通して繰り返すために使用され得る。

【 0 0 5 0 】

他の実施の形態においては、書き込まれているデータのタイプをクライアントプロセス 1 3 4 が特定し得るように、タイプパラメータは、書込コマンド (例えば、write(Key, Value, Context, Type) ResultCode (書込 (キー、値、コンテキスト、タイプ) 結果コード)) に付加され得る。データセットサービス 1 1 2 は、データが最後にアクセスされてから所定量の時間経過後に、(例えば、そのデータが最早必要とされない場合に、記憶領域を再利用するために) そのデータを削除するように設定され得る。削除前に許容される時間は、データのタイプに基づき得る。タイプは、(例えば、データのいくつかのタイ

10

20

30

40

50

プは他のタイプのデータよりも重要であり得るということに基づいて) データセットサービス 1 1 2 が記憶すべきデータのコピーの数を決定するためにも使用され得る。

【 0 0 5 1 】

他の実施の形態において、読出コンテキストは、読出コマンドに対する入力として通過させられることもあり得る(例えば、read(Key, Context) ValueList, Context, ResultCode (読出(キー、コンテキスト) 値リスト、コンテキスト、結果コード))。そのような実施の形態においては、読出コマンドに対する入力として通過した読出コンテキストは、先行する読出の結果として得られることがある。読出動作に対する入力として読出コンテキストを通過させ戻すことにより、クライアントプロセス 1 3 4 は、先行する読出動作の間にアクセスされたデータセットの特定のバージョンの検索において重要性を指示し得る。理解されるであろうように、アプリケーションプログラムインタフェース上における他のバージョンもあり得る。

10

【 0 0 5 2 】

III. データバージョンング (Data Versioning)

A. ベクトルクロックの動作 (Operation of Vector Clocks)

図 1 6 を参照して、データバージョンング構成について説明する。前述したように、高可用性を提供するために、データセットサービス 1 1 2 は、同一データの複数のバージョンが異なるホスト 1 3 0 上に同時に存在することを許容する。データセットの最も新しいバージョンのコピーをその優先リスト 1 9 0 における上位 N のホストに移動するために、進行中の試行が行われるが、このプロセスは瞬間的なものではない。移動が発生する前には、データセットのより古いバージョンのコピーが、その優先リスト 1 9 0 におけるいくつかのホストに、優先リスト 1 9 0 における最上位又は最上位に近いホストにさえ、存在し得る。従って、例えば、一つのホスト 1 3 0 が一時的に失われた古い変更を反映した一つのバージョンを有し、別のホスト 1 3 0 が古い変更が利用不可能になると同時になされた新しい変更を反映した別のバージョンを有することがある。

20

【 0 0 5 3 】

実施の一形態において、同一のデータセットの二つのコピーがそのデータセットの異なるバージョンであって相互に関して差異を有しているか否かを判定することができるようにすることが望ましい。二つのバージョンが相互に先祖と子孫の関係にある(例えば、一方のバージョンが単に古くなっただけで他方のバージョンに組み込まれている)状況を、二つのバージョンがコンフリクト(矛盾)する(例えば、各バージョンが他方のバージョンに反映されていないデータを含んでいる)状況から区別することができるように、それらの差異にアクセスすることができるようにすることも望ましい。

30

【 0 0 5 4 】

実施の一形態において、バージョン履歴は、データセットの各コピーと共に記憶される。例えば、バージョン履歴は、同一データセットの異なるバージョン間の因果関係を捕捉するベクトルクロックの形態で記憶され得る。ベクトルクロックは、二つのバージョンがコンフリクトするか否かを判定することが可能となるように、データセットのバージョン履歴に関する十分な情報を簡潔に記憶し得る。実施の一形態において、ベクトルクロックは、データセットのバージョンに関連する { ホスト ID、カウンタ } ({ host ID, counter }) ペアのリストを含む。ホスト ID (host ID) 値は、書込動作を調整 (コーディネート) したホストを示す。カウンタ (counter) 値は、そのホストがデータセットに書込を行った回数を示している。カウンタ値は、データバージョンについての因果関係、即ち、どのような変更がそのバージョンに先行するかについての要約をコード化 (符号化) する。

40

【 0 0 5 5 】

データセットの二つのバージョンが因果関係を示す順序を有するか(従って一方を無視してもよい)又は並列のブランチ(分岐)であるか(従って調整が必要であるか)の判定を試みる場合、それらのベクトルクロックを検査するだけで十分である。一つのベクトルクロックが、他のベクトルクロックのなかで総てのホスト ID に対して最大の又は最大

50

に等しいカウンタ値を有している場合、前者は後者の子孫であり、後者は無視することができる。従って、ベクトルクロックは、データ展開の複数のブランチをまとめて一つに戻すために、同一データの複数のバージョンを調整することをクライアントプロセス 134 に可能にさせる。

【0056】

図16は、データセットサービス112により使用され得るデータバージョンングの実施例を示している。ステップ400において、最初は、データセットは空(empty)である。ステップ402において、クライアントプロセス134は、ホストAを使用して空データバージョン V_0 を更新する。書込を調整(コーディネート)するホストAは、先行するバージョンのクロックをコピーし、ホストAに関連するカウンタ値を増加させ、データバージョン V_1 に対するベクトルクロックを生成する。この場合、これが最初の更新なので、カウンタは、1にインクリメントされる。データセットサービス112は、データバージョン V_1 及びそれに関連するベクトルクロック $[(A, 1)]$ を記憶する、例えば、ホストAが、ローカル書込動作を実行し、さらに、追加的なローカル書込動作を実行して追加的なコピーを記憶するために新しいバージョンをホストB及びCに送信する。データセットサービス112が商業ウェブサイトシステム内にショッピングカート情報を記憶する例示的な実施の一形態において、この更新は、アイテムをショッピングカートに追加するビジターのために発生し得る。理解されるであろうように、データセットの新しい「バージョン」を構成するものは、アプリケーションに依存して変化し得る。

【0057】

図16において、コーディネータは、優先リスト190における最上位にランクされているNの到達可能なホストの一つである。上述のように、コーディネータは、コーディネータは、優先リスト190における最上位にランクされているNの到達可能なホストの一つではないホスト130であり得る。そのような実施例においては、書込要求を受信したとき、コーディネータは、新しいバージョンに対するベクトルクロックを生成してその新しいバージョンをローカルに記憶するために、そのキーに対する優先リスト190における最上位にランクされているNの到達可能なホストの一つを選択し得る。その後、コーディネータは、既に説明したように、残余の最上位にランクされているNの到達可能なホストにその新しいバージョンを(その新しいベクトルクロックと共に)送信し得る。

【0058】

ステップ404において、同一のクライアントプロセス134は、ホストAを使用してデータバージョン V_1 を更新する。書込を調整(コーディネート)するホストAは、先行するバージョンのクロックをコピーし、ホストAに関連するカウンタ値を2に増加させ、データバージョン V_2 に対するベクトルクロックを生成する。再度、ホストAは、データバージョン V_2 及びその関連するベクトルクロック $[(A, 2)]$ をホストB及びCに転送し、追加的なコピーを記憶する。バージョン V_2 はバージョン V_1 に由来し、従ってバージョン V_1 を上書きするが、まだバージョン V_2 を認識していないホストパーティションに残存しているバージョン V_1 の複製が存在し得る。

【0059】

ステップ406において、要求を調整(コーディネート)するために、同一プロセスがホストBを使用してデータバージョン V_2 を更新する。新しいホストBが更新を調整(コーディネート)するので、新しいベクトルクロックエントリが、このホストBとの関連においてカウンタ値1で生成される。データセットサービス112は、データバージョン V_3 及びその関連するベクトルクロック $[(A, 2); (B, 1)]$ を記憶する。データバージョン V_2 に対するベクトルクロックも、バージョン履歴を保持するために又は実行されるべきより複雑な調整を可能とするために、所望される場合には記憶され得る。ステップ406の後、バージョン V_1 は認識しているがバージョン V_2 は認識していないホストが、バージョン V_3 及びその関連するベクトルクロックを受信することがある。そのホストは、バージョン V_1 及びバージョン V_3 のそれぞれのベクトルクロック $[(A, 1)]$ 及び $[(A, 2); (B, 1)]$ を比較することによって、バージョン V_1 が因果関係に

においてバージョンV₃に先行し従ってバージョンV₃により上書きされるべきことを意味していることを判定することができる。一方、事象の異常なシーケンスが発生して、データバージョンV₃に対するベクトルクロックがバージョンV₁のクロックにおける総てのホストに対して小さい又は等しいカウンタを有する場合、バージョンV₃は、バージョンV₁の先祖であり、抹消され得る。

【0060】

ステップ408において、異なるクライアントプロセス134は、バージョンV₂を読み出し、ホストCを使用してバージョンV₂の更新を試行する。例えば、ホストB及びCは、利用不可能である場合がある。本実施例においては、ホストCが、ステップ406の書込動作に関係しておらず、バージョンV₃を認識していないものとする。新しいホストCが更新を調整(コーディネート)するので、新しいベクトルクロックエントリは、このホストCとの関連においてカウンタ値1で生成される。データセットサービス112は、データバージョンV₄及びその関連するベクトルクロック[(A, 2); (C, 1)]を記憶する。ステップ408において、バージョンV₁又はバージョンV₂を認識しているホストは、バージョンV₄及びその関連するベクトルクロックを受信したときに、バージョンV₁及びバージョンV₂が新しいデータによって上書きされること及び抹消され得ることを、決定し得る。

【0061】

ステップ410において、クライアントプロセス134は、バージョンV₃及びバージョンV₄の両方を読み出す。例えば、読出動作は、ホストAによって調整(コーディネート)されることがあり、かつ、ホストB及びCも関係することがある。ホストAは、ベクトルクロック[(A, 2)]を有するデータセットのそれ自身のコピーと、ベクトルクロック[(A, 2); (B, 1)]を有する、ホストBからのデータセットのコピーと、ベクトルクロック[(A, 2); (C, 1)]を有する、ホストCからのデータセットのコピーとを得る。当該読出のコンテキストは、バージョンV₃及びバージョンV₄のクロックの要約、即ち、[(A, 2); (B, 1); (C, 1)]である。ベクトルクロックの検査からは、バージョンV₃及びバージョンV₄のそれぞれには、相互に反映されていない変更が含まれているので、ホストAは、バージョンV₃とバージョンV₄との間にいかなる因果関係も存在しないことを見出すことになる。従って、バージョンV₃及びバージョンV₄は調整される。

【0062】

実施の一形態において、データセットサービス112(この例では、ホストA)は、調整をどのように実行するかを順番に決定するクライアントプロセス134(及び/又はクライアントプロセス134に関連するバージョン調整ロジック136)に複数のバージョンを提供する。この構成は、調整を実行するために使用される任意のビジネスロジックに、記憶され、又は、データセットサービス112よりむしろクライアントプロセス134と関連付けられることを可能とする。クライアントプロセス134とバージョン調整ロジック136とは分離したものとして示されているが、クライアントプロセス134とバージョン調整ロジック136とが統合された形態で提供され得ることは理解されるであろう。他の実施の形態においては、バージョン調整ロジック136は、データセットサービス112と共に提供され得る。複数のバージョンは、例えば、いずれを保持するかを決定するためのバージョン上のデフォルト序列を使用することにより、単一の調整されたバージョンを生成するために異なる複数のバージョンを併合することにより、データの分析を実行して矛盾をどのように取り扱うかを矛盾ごとに決定することにより、等々により、調整され得る。理解されるであろうように、アプリケーションに依存して、異なる方策が、異なる状況ではより最適であり得る。

【0063】

ステップ412において、書込要求がクライアントプロセス134から受信される。ホストAは、書込を調整(コーディネート)し、ベクトルクロックにおける対応するカウンタ値を更新する。更新されたバージョンは、調整動作に関係しない、クライアントプロセ

10

20

30

40

50

ス 1 3 4 により実行された他の変更を含み得る。新しいバージョン V_5 は、ベクトルクロック $[(A, 3); (B, 1); (C, 1)]$ を有することになる。

【 0 0 6 4 】

データセットに対して調整に加えてどのような変更が実行されたかに拘わらず、ステップ 4 1 2 において、ホスト A がカウンタ値を $[(A, 3); (B, 1); (C, 1)]$ に更新することは理解されるであろう。ベクトルクロック $[(A, 2); (B, 1); (C, 1)]$ を有するいかなるバージョンも存在しないので、ベクトルクロック内のカウンタの更新は、親クロックを新しいクロックから区別する。加えて、複数のクライアントプロセスが同時に調整を試行し得る（例えば、調整のために異なるホストを使用して）が異なる結果に到達し得る（例えば、変更のみならず調整も加える異なる併合ロジック等のために、）ので、カウンタを増加させることが望ましい。もしカウンタが更新されなければ、異なる併合試行が同一クロック即ち $[(A, 2); (B, 1); (C, 1)]$ に割り当てられることがあり、従って、相互に区別することができなくなる。

【 0 0 6 5 】

B . ベクトルクロック情報及び切捨 (Vector Clock Information and Truncation)

実施の一形態において、{ ホスト ID、カウンタ } ({ host ID, counter }) ペアのみを含むよりもむしろ、ベクトルクロックは、多数の追加的な値を含み、以下のような形態を有する。

【 0 0 6 6 】

ベクトルクロック

= { (< ホスト ID > < ホスト情報 > < キー情報 >), < カウンタ > , < タイムスタンプ > }
(VectorClock

= { (< HostID > < host-gen > < key-gen >), < counter > , < time-stamp > })

ホスト ID (host ID) は、ホストに対する固有の識別子であり、カウンタパラメータは、データバージョンに対応する因果関係情報をコード化 (符号化) し、上述した { ホスト ID、カウンタ } ({ host ID, counter }) ペアに対応する。実施の一形態において、(< ホスト ID > < ホスト情報 > < キー情報 >) ((< HostID > < host-gen > < key-gen >)) パラメータの組合せは、ホスト ID 単独に関して前述したような態様で機能する。即ち、三つのパラメータ (< ホスト ID > < ホスト情報 > < キー情報 >) のいずれか一つでも異なっていれば、ホストは、異なるホストであると考えられる (即ち、データセットの異なるバージョン間にいかなる因果関係も伴い得ない) 。

【 0 0 6 7 】

実施の一形態において、ホスト 1 3 0 は、ベクトルクロックを同時にディスクには書き込まない。従って、ホストが各キーに対して生成されたシーケンス番号を見落とすことがあり、その結果そのシーケンス番号が再使用され、それによりベクトルクロックの一貫性が損なわれるという可能性が存在する。見落としの危険性 (例えば、ホスト障害後に) が認識されると、ホスト 1 3 0 は、その < ホスト情報 > (< host-gen >) パラメータを更新し、総ての未来のベクトルクロックに対して、完全に異なるホストとなるように生成し (任意のキーに対して) 、完全に異なるホストとなって見えるようにする。従って、ホスト 1 3 0 が再起動されたときに < ホスト情報 > パラメータをインクリメントすることは、障害の前に生成されたベクトルクロックを、再起動後に生成されたベクトルクロックから区別することを可能にする。理解されるであろうように、各ベクトルクロックに対するカウンタは、無限に単調増加する。実施の一形態において、無限のカウンタ番号を回避するために、各ホストは、例えば、< ホスト情報 > パラメータをインクリメントすることにより、新しい固有の識別性を強制的に、定期的を選択する。例えば、再起動後に新しい固有の識別性を割り当てられたホストは、それにより、< カウンタ > (< counter >) パラメータをゼロに設定もする。これは、識別性を変更する前に単一のホスト 1 3 0 が調整 (コーディネート) できる書込の回数によって、可能性のある最高のカウンタ値が制限される結果をもたらす。他の実施の形態において、識別性の変更は、そのカウンタ値の一つ又は複数

10

20

30

40

50

【 0 0 6 8 】

<キー情報> (<key-gen>) パラメータは、キー生成カウンタをトラック (追跡) するために使用され得る。実施の一形態において、データハンドオフ後に、ホスト 130 は、得られたいかなるデータも削除する。これが、優先リスト 190 において下位にあるホスト 130 のための記憶容量を節約する。同時に、ホスト 130 は、データハンドオフ後にインクリメントされた<キー情報> パラメータを保持し、それにより、次回にホスト 130 が書込動作を実行することを要求されたと仮定した場合におけるいかなる因果関係も回避する。例えば、ホスト D が、ベクトルクロック [(A, 3), (D, 1)] を有するデータセットのバージョンに対する書込動作を調整 (コーディネート) し、データハンドオフを実行し、その後、ベクトルクロック [(A, 2)] を有するデータセットのバージョンに対する別の書込動作を調整することを要求される場合、更新されたデータセットがベクトルクロック [(A, 3), (D, 2)] を有するようにすることは不適當である。この状況において新しい<キー情報> 値を割り当てることにより、ホスト 130 は、新しいホストであるかのごとく見えるようになり、それにより、二つのバージョン間の因果関係の出現を回避する。実施の一形態において、各ホスト 130 は、キーごとに別個の<キー情報> を保持し、また、ベクトルクロックは変更された識別性 (例えば、変更された<ホスト ID> 又は更新されたその<ホスト情報>) を保持するので、各ホスト 130 は、対応してベクトルクロックが生成された総てのキーについてキーの世代を記憶する。同様に、各ホスト 130 は、対応する<キー情報> パラメータ又は<ホスト情報> パラメータのいずれかが更新されてから、キーに対するベクトルクロックにおいて使用された最後の<カウンタ> パラメータも記憶し得る。

10

20

【 0 0 6 9 】

<タイムスタンプ> パラメータは、データセットの世代及びそのベクトルクロック内のエントリを監視するために使用され得る。いくつかのアプリケーションにおいては、データが所定の世代を超えた場合、そのデータを削除することが望ましい。例えば、ショッピングカートアプリケーションにおいて、数日、数週間、数月又は数年等の期間に亘って放置されているショッピングカートは、削除することが望ましい場合がある。タイムスタンプは、この方法におけるデータセットの削除を支援 (サポート) するために使用され得る。加えて、タイムスタンプは、ベクトルクロック切捨 (Vector Clock Truncation) のためにも使用され得る。理解されるであろうように、データセットとの関係において書込動作が調整 (コーディネート) された異なるホスト (又は異なる<ホスト情報> 若しくは<キー情報> を有する同一ホスト) のリストの長さが増加するのに伴い、そのデータセットに対するベクトルクロックの長さも増加する (即ち、ベクトルクロックに含まれる {ホスト ID、カウンタ} ({host ID, counter}) ペアのリストの長さが増加するからである)。従って、所定数だけ世代を重ねたベクトルクロックは、タイムスタンプを使用して、削除され又は切り捨てられ得る。

30

【 0 0 7 0 】

他の実施の形態においては、データセットにおける変更をトラック (追跡) するために、ベクトルクロックを使用するよりむしろ、他のバージョン履歴機構が使用され得る。例えば、ハッシュ履歴も使用され得る。ここで、用語「バージョン履歴 (version history)」とは、データセットにおける変更を経時的にトラック (追跡) するために (即ち、変更が存在することをトラックするために、であって、必ずしも変更の性質をトラックするために、ではない) 使用され得る任意のデータ構造をいう。理解され得るように、異なるバージョン履歴機構は、ディスク領域使用量、バンド幅、古いバージョンが削除されたときの一貫性の保持、因果関係における先行の検出の速度及び容易さ、等による異なるトレードオフを提供し得る。実施の一形態において、データセットの二つ又はそれ以上のコピー間の因果関係における先行 (又はその不存在、上記ではコンフリクト (矛盾) と称された) の検出を可能とするバージョン履歴機構が使用される。バージョン履歴機構は、バージョンコンフリクトを許容して、データの損失なしに (可用性) を発生させるために、及び、優先リストにおける上位のホストへのデータ移動の際に一貫性の保持を容易にするた

40

50

めに、使用され得る。

【 0 0 7 1 】

I V . 複数データセンタ (Multiple Data Centers)

A . 複数データセンタ構成の構造 (Architecture of Multiple Data Center Arrangement)

図 1 7 乃至図 2 5 を参照すると、データ処理システム 1 0 0 の他の実施の形態が示されている。図 1 7 乃至図 2 5 において、データセットは、複数段 (multi-tiered) リング構成に従ってホスト 1 3 0 間に分割されている。複数段リング構成は、例えば、各ホストが異なる地理的ロケーションに (例えば、異なる都市、異なる国、異なる大陸に存在し得る異なるデータセンタに) 配置されているデータセット記憶システムを実施するために使用され得る。例えば、ホスト間における相互関係障害の確率を低減するために、データは、そのような異なるデータセンタ間で複製され得る。単一のデータセンタの障害が全システムの可用性に重大な影響を与える可能性は非常に小さい。加えて、クライアント要求をより近いデータセンタに転送することにより (ネットワーク待ち時間 (network latency) によって)、端末間 (end-to-end) データ検索応答時間は短縮され得る。複数段リング構成は、他の理由のために、例えば、共有データセンタ内に位置するホスト等においても使用され得る。例えば、リングの異なる段は、データセンタ内の領域、データセンタ内のホストの特定のラック (rack) 等を特定するために使用され得る。一つの実施例を提供する目的のために、図 1 7 においては、各ホストが異なるデータセンタに配置されているデータセット記憶システムを実施するために、複数段リング構成が使用されることが想定されている。

【 0 0 7 2 】

先ず図 1 7 を参照すると、図 1 7 は、データ処理システム 1 0 0 が 2 段又は 2 レベルリング構成を備えている実施の形態を示している。2 段リング構成は、図 1 及び / 又は図 2 に示されたデータセットサービス 1 1 2 を実施するために使用され得る。図 1 7 において、データ処理システム 1 0 0 は、上位レベルリング 5 0 4 上に論理的に配置されている複数のデータセンタ 5 0 2 を備えている。四つのデータセンタ 5 0 2 が示されているが、実際には、任意の数のデータセンタ 5 0 2 が使用され得ることは理解されるであろう。

【 0 0 7 3 】

データセンタ 5 0 2 は、通信ネットワーク 5 0 8 によって相互に接続され得る (例えば、広域ネットワーク、インターネット等)。データセンタ 5 0 2 間のメッセージ通信は、図 2 5 との関連において詳細に後述されるように、メッセージフィルタ 5 1 0 を通過し得る。図 1 乃至図 2 におけるように、データセンタ 5 0 2 のそれぞれは、通信ネットワーク 1 0 4 (例えば、インターネット) 経由で種々のユーザコンピュータ 1 0 2 によってアクセスされ得る。

【 0 0 7 4 】

データセンタ 5 0 2 のそれぞれは、各下位レベルリング 1 8 4 上に論理的に配置されている複数のホスト 1 3 0 をさらに備えている。図示されている実施例においては、各下位レベルリング 1 8 4 は、異なるデータセンタ 5 0 2 に対応している。各データセンタ 5 0 2 内では、各リング 1 8 4 上のホスト 1 3 0 はまた、図 3 乃至図 1 6 との関連において上述したように動作し得る。下位レベルリング 1 8 4 は、同種又は異種 (例えば、異なる数のホスト、異なるハッシュ関数、異なる構成等を有する) であり得る。さらに、以下に説明するように、データセンタ 5 0 2 に関する上位レベルリング 5 0 4 の動作は、ホスト 1 3 0 に関して図 3 乃至図 1 6 との関係において上述した通りのリング 1 8 4 の動作と同一であり得る。

【 0 0 7 5 】

図 1 8 を参照すると、実施の一形態において、データセット記憶システム 1 1 8 は、各データセンタ 5 0 2 がデータセットの一部を記憶するように、各データセンタ 5 0 2 に分散され得る。データセンタ 5 0 2 のそれぞれは、上位レベルリング 5 0 4 上におけるハッシュ値の一つの範囲 (又は、詳細に後述するように、上位レベルリング 5 0 4 上における

ハッシュ値の範囲の組（複数の範囲））に対して責任を有するものとすることができ、その場合、ホスト130及びリング184との関連において上述したと同様の態様で、各データセンタ502は、ハッシュ範囲におけるそれ自体の位置から先行するデータセンタ502の位置までに延在するハッシュ値との関係において読出／書込動作に責任を有する。データセットへのアクセスの要求が受信されたとき（例えば、読出動作又は書込動作を通じて）、データセットがアクセスされ得るデータセンタ502を決定するために、上位レベルリング504のためのハッシュ関数にキーが適用される。（図18において、符号DC1-DC4はそれぞれ、図17における四つのデータセンタ502の異なる一つを示している。）データセットがアクセスされ得るホスト130を、関係するデータセンタ502内で決定するために、キーは、下位レベルリング184のためのハッシュ関数にも適用される。上位レベルリング504のために使用されるハッシュ関数は、下位レベルリング184のために使用されるハッシュ関数と同一のもの又は異なるものであり得る。同様に、上述のように、下位レベルリング184のそれぞれのために使用されるハッシュ関数は、他の下位レベルリング184のために使用されるハッシュ関数と同一のもの又は異なるものであり得る。図18に示された構成のマッピングでは、個々のデータセンタ502は、データセンタ502へのデータセットの分割についての全体的な再マッピングを行うことなく、追加し又は除去することが可能であり、それによりスケーラビリティ（拡大縮小可能性）（scalability）が増進される。

10

【0076】

実施の一形態において、データセンタ間のデータ複製も、ホスト130との関連において図10について上述したと同様の態様により、支援（サポート）され得る。従って、図18に示したように、リング504上の直近のデータセンタ502に単純にデータセットが割り当てられるよりむしろ、データセットは、最初のMの後続するデータセンタ502に割り当てられるものとしてもよい。データセットサービス112は、データセットがMのデータセンタにおいて複製されることを保証するために動作することが可能であり、各データセンタ502は、それ自体とそのM番目の先行するデータセンタ502との間のリング504上の範囲に対して責任を有し得る。

20

【0077】

所定のデータセットの複製を記憶するデータセンタ502の数は、例えば、データセットごとの基準、データタイプごとの基準等で設定可能であるものとし得る。理解されるであろうように、保持される各データセットの複製の数は、特に、可用性の望ましいレベル及び通信ネットワーク508上の更新通信量の望ましいレベルに基づいて決定され得る。即ち、異なるデータセンタに亘ってより多くの複製が記憶されるほど、可用性は増加する。しかし、データセットの複製されたコピーを、一貫性を有するように保持するための更新の間に、通信ネットワーク508上のネットワーク通信量も増加する。データセットが一つのデータセンタ502内で複製されるべきものと仮定すると、データセットを複製するデータセンタ502内のホストの数も、例えば、データセンタごとの基準、データセットごとの基準、データタイプごとの基準等に基づいて、設定可能であるものとし得る。

30

【0078】

実施の一形態において、データセンタ間の負荷平衡も、ホスト130との関連において図12について上述したと同様の態様により、支援（サポート）され得る。例えば、データセンタ502は、リング504上の複数の位置に割り当てられ得る。そのような構成は、当該構成を使用しなかったとするとリング504上における各データセンタ502の無作為の位置割り当てによって発生し得る不均一なデータ及び負荷の分散を回避するために使用され得る。そのような複数配置は、各データセンタ502に割り当てられるデータセットの数の不均一を低減するように作用する。その理由は、リング184上における無作為配置の増加した数は、各データセンタ502に割り当てられたデータセットの数を平均値に収束させるように作用するからである。加えて、リング504上の複数の位置へのデータセンタ502の割当ては、異種のデータセンタの使用を容易にもする、即ち、より性能の高いデータセンタ502（例えば、処理能力、記憶容量及び／又はネットワーク容

40

50

量に基づいて決定される)がリング504上のより多くの位置を割り当てられ、より性能の低いデータセンタ502がリング504上のより少ない位置を割り当てられるようにし得る。さらに加えて、リング504上の複数の位置へのデータセンタ502の割当ては、データセンタ間での負荷の移動も容易にする。その理由は、各データセンタ502が、他のデータセンタ502のそれぞれと後続/先行関係を有し得るからである(リング504上の各データセンタ502に十分な数の位置が割り当てられているものとする)。従って、例えば、データセンタ502のうちの一つが利用不可能又は非動作状態になったとすると、非動作状態になったデータセンタ502により処理されていた負荷は、データ可用性を損なうことなく、残余の利用可能な各データセンタ502にほぼ均等に分散させられ得る。

10

【0079】

図19を参照すると、各データセットは、各データセンタ502がキーに基づいて生成されたハッシュ値からリング504上を時計回りに周回するときに最初に遭遇する他のデータセンタ502の優先リスト519を有し得る。優先リスト519は、データセットへのアクセス(例えば、読出、書込等)に使用されるデータセンタ502の好適な順序を表している。総てのデータセンタ502が利用可能である場合、優先リスト519における上位Mのデータセンタ502がデータセットを記憶する。同一データセットにおける連続する動作がMのデータセンタの同一セットにアクセスすることがあり、従って一貫性を有するものであり得る(即ち、同一キーにおける先行する動作によって読出/書込が行われた同一データに、動作がアクセスする)。優先リスト519における一つ又は複数のデータセンタに故障が発生した場合、又は、ネットワークパーティションが存在する場合、データセットは、優先リスト519における下位にランク付けされた一つ又は複数のデータセンタ502に一時的に記憶され、それにより高可用性が維持される。加えて、同一データセットへの連続する動作がデータセンタ502の異なるセットにアクセスすることがあるが、アクセスされるデータセンタ502のセットにいくらかの重複がある限り、動作は依然として一貫性を有し得る。優先リスト519において上位の利用可能なデータセンタ502にアクセスすることにより、動作から動作へのホストの利用可能性における小さい変化は、一貫性に否定的な影響を与えない。その理由は、後続のアクセスが、重複するデータセンタを含み得るからである。

20

【0080】

優先リスト519は、例えば、ハッシュ関数に基づいて算定され得る。実施の一形態において、所定のデータセットを記憶しているデータセンタ502にアクセスするために、各ホスト130は、データセンタ502の利用可能性(優先リスト519における上位Mの利用可能なデータセンタを選択するために)に加えて、ハッシュ空間におけるデータセンタ位置に関する情報(優先リスト519を算定するために)も記憶し得る。他の実施の形態においては、例えば、記憶される優先リスト519が、ハッシュ関数に基づいて構成されること、及び、優先リスト519を構成する際に考慮に入れることが望ましい場合がある他の関数に基づいて構成されることを可能とするために、優先リスト519は、記憶され得る。

30

【0081】

B. アクセス動作(Access Operations)

図20乃至図24を参照すると、データセンタ502に記憶されているデータセットへのアクセスに関連する動作が示されている。図20は、実施の一形態に係る図17のシステムにより実行されるアクセス動作のフローチャートである。図21乃至図24は、実施の一形態に係る図20のアクセス動作の態様を詳細に示す図である。

40

【0082】

ステップ602において、ユーザコンピュータ102との接続がデータセンタ502により確立される。理解されるであろうように、各データセンタ502は、データセットサービス112を実現するホスト130だけでなく、ネットワークインタフェース110及び他のサービス114を実現する他のホストも含み得る。従って、図21を参照すると、

50

例えば、ネットワークインタフェース 1 1 0 を実現するホストの一つであり得るホスト 5 3 2 との接続が確立され得る。

【 0 0 8 3 】

実施の一形態において、ユーザコンピュータ 1 0 2 との接続は、一つのデータセンタ 5 0 2 (例えば、無作為基準による可能性もある)において確立される場合があり、その後、他のデータセンタ 5 0 2 に転送される。例えば、図 2 1 において、ユーザコンピュータ 1 0 2 との接続は、一つのデータセンタ D C 4 内のホスト 5 3 2 によって確立されることがあり(ステップ 6 0 2)、その後、例えば、より近いものであってもよく、より少ない負荷しか掛かっていないものであってもよく、及び/又は、接続を維持するためにより適したものとなるような他の特性を示すものであってもよい、他のデータセンタ D C 1 内の他のホスト 5 3 4 に転送される(ステップ 6 0 4)。

10

【 0 0 8 4 】

ステップ 6 0 6 において、データアクセス要求(例えば、読出要求、書込要求等)が受信される。図 2 2 を参照すると、データアクセス要求は、クライアントプロセス 1 3 4 を実行しているものであり得るホスト 5 3 6 から、データセットサービス 1 1 2 内のホスト 1 3 0 によって受信され得る(図 3 参照)。例えば、上述の図 2 において示した実施例との関係においては、ホスト 5 3 4 は、ネットワークインタフェース 1 1 0 を実現するホストの一つである場合があり、ユーザコンピュータ 1 0 2 に接続されている場合があり、ホスト 5 3 6 は、ショッピングカートサービス 1 2 4 を実現するホストの一つである場合があり、ホスト 5 3 4 から要求を受信する場合があり、ホスト 1 3 0 は、データセットサービス 1 1 2 を実現するホストの一つである場合があり、ホスト 5 3 6 からアクセス要求を受信する場合がある。データセットへのアクセス要求がデータセンタ 5 0 2 内のホスト 1 3 0 において受信されると、ステップ 6 0 8 において、ホスト 1 3 0 は、データセットがデータセンタ 5 0 2 内にローカルに記憶されているか否かを判定する。データセットは、データセンタ 5 0 2 内にローカルに記憶され得る。その理由は、例えば、データセンタ 5 0 2 が優先リスト 5 1 9 における上位 M のデータセンタの一つだからであり、データセンタ 5 0 2 が優先リスト 5 1 9 における下位のデータセンタであるが優先リスト 5 1 9 における上位 M のデータセンタの一つであるデータセンタ 5 0 2 にデータセットを移動するまで一時的にデータセットを記憶しているからであり、データセンタ 5 0 2 がユーザとの接続を確立してデータセットの貸し出された(leased)コピーを一時的に記憶しているからであり(詳細に後述するように)、又は、別の理由があるからである。データセットがローカルに記憶されている場合は、ステップ 6 1 0 において、応答は、データセットの一つ又は複数のローカルコピー(場合によっては一つ以上のバージョン)に基づいて供給され得る。そうでない場合は、ホスト 1 3 0 は、他のデータセンタ 5 0 2 からデータセットの一つ又は複数のコピー(場合によっては一つ以上のバージョン)を取得し得る。データセットのコンフリクト(矛盾)するバージョン(例えば、一つのデータセンタ内からのコンフリクトバージョン、異なるデータセンタからのコンフリクトバージョン、又は、その両方)が存在する場合、そのようないかなるコンフリクトバージョンも、データセットを要求する特定のクライアントプロセス 1 3 4 に関連するデータセットバージョン調整ロジック 1 3 6 に報告され、上述のように、データセットバージョン調整ロジック 1 3 6 により解決され得る。例を挙げる目的のために、データセンタ D C 1 が、データセットの複製を記憶する M のデータセンタの一つではないものと仮定する。従って、データセンタ D C 1 におけるホスト 1 3 0 が、他のデータセンタからデータセットのコピーを得るためにコーディネータとして動作する。

20

30

40

【 0 0 8 5 】

ステップ 6 1 2 において、アクセス要求が受信された後、データセットに対するキーが、上位レベルリング 5 0 4 及び下位レベルリング 1 8 4 に対するハッシュ関数に適用される。ステップ 6 1 4 において、図 2 2 を参照すると、データセンタ D C 1 におけるホスト 1 3 0 (コーディネータとして動作する)が、優先リスト 5 1 9 における一つ又は複数の上位のデータセンタからのデータを要求する。実施の一形態において、ホスト 1 3 0 は、

50

上位レベルリング504に対するハッシュ関数にキーを適用して、データセンタDC2及びDC3にアクセス要求を送信する(例えば、データセンタDC2及びDC3がそのデータセットに対する優先リスト519の上位にあることが判定された後に)。アクセス要求がデータセンタDC2及びDC3における各ホスト130により受信されると、それらのホスト130は、下位レベルリング184に対するハッシュ関数にキーを適用して、そのデータセットを記憶する各データセンタ内のホスト130を判別する。この手法において、データセンタDC1におけるホスト130にとって、遠隔のデータセンタDC2及びDC3のリング184上におけるホスト130の位置に関する情報を記憶することは、必要ではない。他の実施の形態においては、各データセンタ502における各ホスト130はこの情報を記憶し、データセンタDC1におけるホスト130は、上位レベルリング504及び下位レベルリング184の両方に対してキーを適用し得る。

10

【0086】

実施の一形態において、ユーザコンピュータ102との接続が確立されたときに、データセットはプリフェッチ(pre-fetch:先取り)され得る。例えば、図2のショッピングカートの実施例との関係においては、ユーザコンピュータ102は、ホスト534との接続を確立し得るが、それは、ショッピングカートデータセットに対する要求が行われるいくらか前であり得る。例えば、ユーザは、ショッピングカートデータセットに対してアクセスすることを必要とするアクション(動作)が実行される前の時間は買い物を行い得る。従って、ショッピングカートデータセットに対してアクセスすることを必要とするアクションをユーザが実行するのを待たずに、ユーザコンピュータ102との接続が確立されると直ちに、データセットは、遠隔のデータセンタ502からの読出動作を実行することによりプリフェッチされ得る。この構成は、通信ネットワーク508経由でデータセット取得することに関連するネットワーク待ち時間(network latency)を回避するために用いられ得る。

20

【0087】

ステップ616において、遠隔のデータセンタ502は、アクセス要求を処理し、データセンタDC1におけるホスト130により受信される応答を送信する。実施の一形態においては、良好な読出動作のために、読出動作は、 R_{DC} のデータセンタにおいて良好に行われなければならない。ここで、 R_{DC} は、設定可能な値であって、 $R_{DC} \leq M$ である。実施の一形態においては、良好と考えられる書込動作のために、書込動作は、 W_{DC} のデータセンタにおいて良好に行われなければならない。ここで、 W_{DC} は、設定可能な値であって、 $W_{DC} \leq M$ である。 $R_{DC} + W_{DC} > M$ となるような R_{DC} 及び W_{DC} の設定は、読出及び書込動作に関係するデータセンタ502のセットの間における重複の設定可能な高い確率が存在するクォーラムのようなシステム(quorum-like system)をもたらす。

30

【0088】

理解されるであろうように、データセンタ502からのデータセットにアクセスする場合、ホスト130に関して上述したように、データセットは、データセンタ502の同一のセットに書き込まれる必要はなく、また、データセンタ502の同一のセットから読み出される必要もない。例えば、データセットは、優先リスト519における下位のデータセンタ502に書き込まれて、優先リスト519における上位のデータセンタ502にデータハンドオフを通じて移動され、その後、最終的に優先リスト519における上位のデータセンタ502から読み出され得る。この態様において、優先リスト519の上位Mのデータセンタにおけるデータセットの最終的な一貫性は、実現され得る。データセンタ502はまた、それらが共有している範囲の下位レベルデータベース比較を周期的に実行し、その後、比較の間に検出されたいかなる差異も調整するために必要なデータ転送を実行する(例えば、データセットの喪失したコピーのために)。従って、データセットサービス112は、それらの優先リスト519における上位Mのデータセンタにデータセットの最新バージョンのコピーを動的に移動させる進行中の試行を行い得る。たとえデータセットの最新バージョンのコピーが、その優先リスト519の下位にあるデータセンタ502

40

50

に最初はコピーされることがあり、又は、別の理由により上位Mのデータセンタの一つにおいて喪失することがあるとしても、そのコピーは、優先リスト519における上位Mのデータセンタへ最終的に移動し戻され、上位Mのデータセンタにおけるデータセットの最終的な一貫性に帰着する。

【0089】

ステップ618において、種々のホスト130及びデータセンタ502から受信される総てのデータセットについてのバージョン履歴は、異なるデータセンタから受信されるデータセット間の一貫性をチェックするために比較される。実施の一形態において、バージョン履歴はベクトルクロックであり、図16との関連において上述したベクトルクロック構成は、異なるデータセンタ502に記憶されている同一データセットの異なるバージョン間の因果関係を捕捉するために使用される。例えば、いかなる所与の二つのホスト130も、それらが異なるセンタ内のものであったとしても、相互に区別され得るように、総てのホスト130は、普遍的に固有の<ホストID>(<Host ID>)を与えられ得る。そのような構成においては、データバージョンングを実行するために使用されるロジックは、ホスト130が複数段(multi-tiered)リング構成に従って組織的構造を与えられているという事実を必ずしも認識している(又は考慮に入れている)必要はない。他の実施の形態において、データバージョンングは、下位レベルリング184のレベルにおいて、及び、上位レベルリング504のレベルにおいて、分離して実行される。そのような実施の形態においては、ベクトルクロックは、書込動作をコーディネート(調整)したデータセンタを示す<データセンタID>(<data center ID>)パラメータを含み得る。

【0090】

ステップ620において、図24を参照すると、データセットは、遠隔のデータセンタにおいてリフレッシュされる。実施の一形態において、一旦データセットが取得されると、データセットは、ある時間期間(ここでは「貸出時間(lease time)」と称する)に亘ってデータセンタDC1に保持される。その後、貸出(lease)が未だ終了していないという条件で、未来の読出動作がローカルに実行される。例えば、データセンタがデータセットへの読出を受信すると、そのデータセンタは、読出要求を送信することなく、ローカルホスト130からのデータセットを遠隔のデータセンタ502へ返送する。書込動作に関しては、実施の一形態において、書込動作は、上述したように、メッセージフィルタを使用して制御され得る。他の実施の形態において、書込動作は、読出動作について上述したのと同様の貸出構成を使用して実行され得る。即ち、データセンタがデータセットへの更新を受信すると、ホスト130は、その更新をローカルに実行し、貸出の終了時にのみ、その更新を他のデータセンタ502に非同期的に伝達する。データセットのローカルコピー上で動作を行うことにより、ユーザコンピュータ102において体験される待ち時間は、減少させられる。他の実施の形態においては、メッセージフィルタも貸出構成も使用されず、更新は、直ちに他のデータセンタ502に伝達される。

【0091】

ネットワークサービス、例えば、小売ウェブサイトとの関係において、貸出時間は、ユーザコンピュータ102の平均セッション時間に基づいて決定されてもよく、ユーザコンピュータ102とのセッションが実行中(アクティブ)であるときは延長され得る。しかし、この貸出時間の間は、データセットが複数のデータセンタにおいて同時にアクセスされる場合(例えば、複数のユーザコンピュータ102が、異なるデータセンタにおける同一のデータセットにアクセスしている場合)、コンフリクトがデータセットの異なるコピーにおいて発生する可能性もある。そのような状況においては、ローカルコピーをリフレッシュし(図24に示されるように)、そのローカルコピーを他のデータセンタに保持されている他の複製コピーと同期させることが望ましいことであり得る。異なるデータセンタ502に記憶されているデータセットにコンフリクトが発生するという点については、そのようなコンフリクトは、データセットバージョン調整ロジック136によって解決され得る。

【0092】

C. メッセージフィルタ動作 (Message Filter Operations)

図 25 を参照すると、実施の一形態において、例えば、データセンタにデータを記憶するためのコストがネットワーク通信量（例えば、ピークネットワーク通信量、平均ネットワーク通信量等）のレベルによって影響を受ける場合、可用性の望ましいレベルになお適合させながら、データセンタ間の通信オーバーヘッドを減少させることは望ましいことであり得る。例えば、他のデータセンタへのデータセットの即時の伝達を回避して、通信（トラフィック）バーストが平滑化されるようにすることは、望ましいことであり得る。

【0093】

そのために、データセンタ 502 のそれぞれにおけるメッセージフィルタ 510 は、ネットワーク通信量を変調するために使用され得る。各メッセージフィルタ 510 は、対応するデータセンタ 502 内においてロジック的（論理的）に集中化され得る。メッセージフィルタ 510 は、データセンタ 502 内のホスト 130 からの書込要求を受信して、その書込要求を即時に又は遅延した態様で伝達するように構成され得る。図 25 に示されるように、データセンタ 502 は、バースト 552 を伴うネットワーク通信量を発生する。メッセージフィルタ 510 は、バーストを平滑化して、ネットワーク通信量が時間期間に亘って分散させられた帯域幅波形 554 を生成するために動作する。他の実施例として、メッセージフィルタ 510 は、データセットのより古いバージョンを廃棄して（ベクトルクロックの分析に基づいて）、データセットの最新のバージョンのみを転送するために構成され得る。メッセージフィルタ 510 はまた、信頼性の高いメッセージ記憶システムとしても構成され得る。例えば、データセンタ 502 がダウンした場合又はメッセージを受信するために利用不可能になった場合、メッセージフィルタ 510 は、メッセージを記憶し、データセンタ 502 がオンラインに復帰したときにそのメッセージを送信するように構成され得る。

【0094】

本発明は、図面を参照して、以上に説明されている。それらの図面は、本発明に係るシステム、方法及びプログラムを実施する特定の実施の形態についてのある一定の詳細を示している。しかし、図面による本発明の記載は、図面中に存在し得るいかなる限定も本発明に課するものと解釈されるべきではない。本発明は、方法、システム、及び、その動作を達成するための任意の機械読取可能な記録媒体上のプログラム製品を企図している。本発明の実施の形態は、既存のコンピュータプロセッサを使用して、又は、この若しくは他の目的のために組み込まれた専用コンピュータプロセッサにより、又は、結線接続されたシステムにより、実施され得る。

【0095】

上述したように、本発明の範囲内にある実施の形態は、その上に記録された機械実行可能な命令若しくはデータ構造を運搬又は保持するための機械読取可能な記録媒体を含むプログラム製品を含む。そのような機械読取可能な記録媒体は、汎用若しくは専用のコンピュータ又はプロセッサを有する他の機械によりアクセス可能な任意の利用可能な記録媒体とすることができる。例として、そのような機械読取可能な記録媒体は、RAM、ROM、EPROM、EEPROM、CD-ROM 若しくは他の光学ディスク記録媒体、磁気ディスク記録媒体若しくは他の磁気記録装置、又は、機械実行可能な命令若しくはデータ構造の形態でのプログラムコードを運搬若しくは記憶するために使用可能であり、汎用若しくは専用コンピュータ又はプロセッサを有する他の機械によりアクセス可能な任意の他の記録媒体を含み得る。情報が、ネットワーク若しくは他の通信接続（結線接続、無線のいずれか、又は、結線接続若しくは無線の組合せ）を通じて機械に転送され又は供給されると、その機械は、機械読取可能な記録媒体としての接続を適当に調査する。従って、そのような任意の接続は、機械読取可能な記録媒体と適当に称される。上記の組合せも、機械読取可能な記録媒体の範囲内に含まれる。機械実行可能な命令は、例えば、汎用コンピュータ、専用コンピュータ又は専用処理機械に、ある一定の機能若しくは機能のグループを実行させるための命令及びデータを含む。

【0096】

本発明の実施の形態は、例えば、ネットワーク接続環境における機械によって実行されるプログラムモジュールの形態におけるプログラムコード等の、機械実行可能な命令を含むプログラム製品による実施の一形態において実施され得る方法ステップの一般的な前後関係において記載されている。概して、プログラムモジュールは、特定のタスクを実行し又は特定の抽象データタイプを実施するルーチン、プログラム、オブジェクト、コンポーネント、データ構造等を含む。機械実行可能な命令、関連するデータ構造、及び、プログラムモジュールは、ここに開示されている方法のステップを実行するためのプログラムコードの例を表している。そのような実行可能な命令又は関連するデータ構造の特定のシーケンスは、そのようなステップにおいて説明された機能を実施するための対応する動作の例を示している。

10

【0097】

本発明の実施の形態は、プロセッサを有する一つ又は複数の遠隔コンピュータへの論理的（ロジック的）結合を使用するネットワーク接続環境において実施され得る。論理的結合は、ここでは例として示されるものであって限定ではないローカルエリアネットワーク（LAN）及びワイドエリア（広域）ネットワーク（WAN）を含み得る。そのようなネットワーク接続環境は、職場全域若しくは事業全域コンピュータネットワーク、イントラネット及びインタネットにおいて一般的であり、異なる通信プロトコルの幅広い変形を使用し得る。当該技術分野における通常の知識を有する者は、そのようなネットワーク計算方式環境が典型的には、パーソナルコンピュータ、ハンドヘルド（手持ち式）装置、マルチプロセッサシステム、マイクロプロセッサベースの又はプログラム可能なコンシューマ（消費者）電子機器、ネットワークPC、サーバ、ミニ（小型）コンピュータ、メインフレームコンピュータ等を含む多くの種類のコンピュータシステム構成を包含することを理解するであろう。従って、図1に図示されたユーザコンピュータ102は、デスクトップコンピュータ、ラップトップコンピュータ、セットトップボックス、携帯（個人）情報端末（personal digital assistant：PDA）、携帯電話、メディアプレーヤ、ウェブパッド、タブレット等を含み得るが、それらには限定されない。本発明の実施の形態は、通信ネットワークを通じて（結線接続リンク、無線リンクのいずれかにより、又は、結線接続リンク若しくは無線リンクの組合せにより）連結（リンク）されている局所（ローカル）及び遠隔（リモート）の処理装置によりタスクが実行される分散型計算方式環境においても実施され得る。分散型計算方式環境においては、プログラムモジュールは、局所（ローカル）及び遠隔（リモート）のメモリ記憶装置に配置され得る。

20

30

【0098】

本発明の全体的なシステム又は部分を実施するための例示的なシステムは、演算処理装置、システムメモリ、及び、当該システムメモリを含む多様なシステム構成要素（コンポーネント）を当該演算処理装置に結合するシステムバスを含む、コンピュータの形態における汎用計算方式装置を含み得る。システムメモリは、読出専用メモリ（ROM）及びランダムアクセスメモリ（RAM）を含み得る。コンピュータは、磁気ハードディスクからの読出及び磁気ハードディスクへの書込を行うための磁気ハードディスク駆動装置、取り外し可能磁気ディスクからの読出及び取り外し可能磁気ディスクへの書込を行うための磁気ディスク駆動装置、CD-ROM又は他の光学記録媒体等の取り外し可能光学ディスクからの読出及び取り外し可能光学ディスクへの書込を行うための光学ディスク駆動装置を含み得る。駆動装置及びそれらの関連する機械読出可能な記録媒体は、機械実行可能な命令、データ構造、プログラムモジュール、及び、コンピュータのための他のデータの揮発性記憶装置を提供する。

40

【0099】

ここで提供されているフローチャートは方法ステップの特定の順序を示しているが、それらのステップの順序は図示されているものとは異なり得ることが理解される点に留意すべきである。また二つ又はそれ以上のステップは、同時に又は部分的に同時に実行され得る。そのような変更は、選択されたソフトウェア及びハードウェアシステム並びに設計者の選択に依存することになる。総てのそのような変更は、本発明の範囲内にあることが理

50

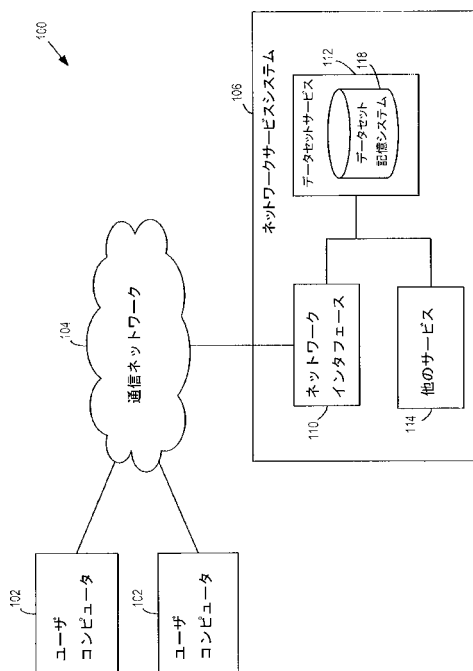
解される。同様に、本発明のソフトウェア及びウェブ手段は、種々のデータベース検索ステップ、相関ステップ、比較ステップ及び決定ステップを達成するための、規則に基づくロジック及び他のロジックによる標準プログラミング技術によって達成され得る。ここで及び特許請求の範囲で使用される単語「エンジン(engine)」は、ソフトウェアコードの一つ又は複数のラインを使用する手段、及び／又は、ハードウェア手段、及び／又は、手動入力を受信するための装置を包含することが意図されているということも留意されるべきである。エンジン、インタフェース、データベース、ブラウザ等の構成要素(コンポーネント)は、相互に通信し得る。そのような構成要素は、統合された態様で提供されるからであり、それらは、ネットワーク等の通信リンクを通じて相互に通信するからであり、及び／又は、その他の理由によるからである。

10

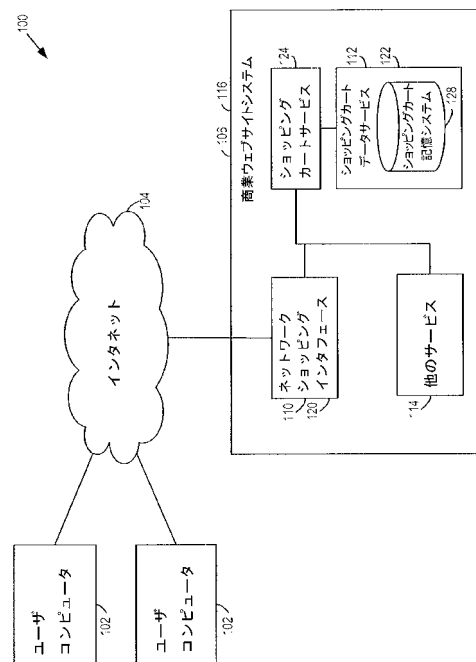
【0100】

本発明の実施の形態についての以上の記載は、例示及び説明の目的のために提示されているものである。総てを網羅すること又は本発明を開示された通りの正確な形態に限定することは意図されておらず、変形及び変更は、上記教示に鑑みて可能であり、本発明の実施から獲得され得る。実施の形態は、当該技術分野における通常の知識を有する者が本発明を企図された特定用途に適当な種々の態様で及び種々の変形と共に利用することができるように、本発明の原理及びその実施上の応用を説明するために選択されて記載されたものである。

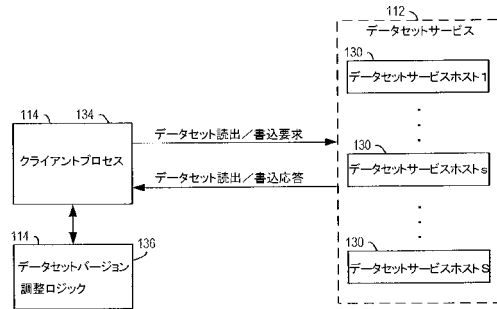
【図1】



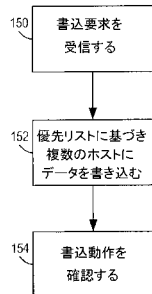
【図2】



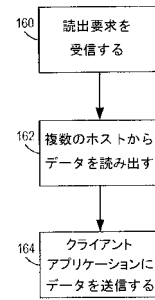
【図 3】



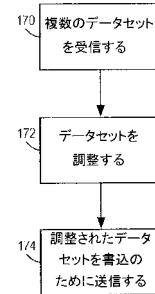
【図 4】



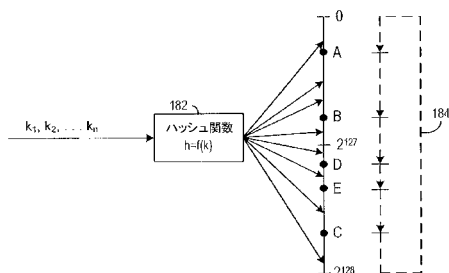
【図 5】



【図 6】



【図 7】



【図 9】

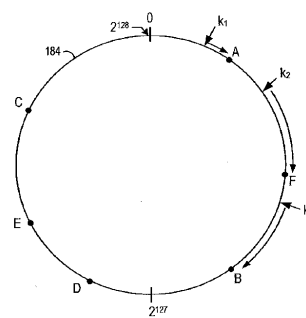


FIG. 9

【図 8】

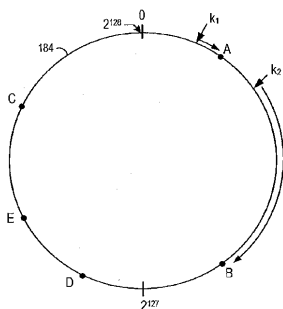


FIG. 8

【図 10】

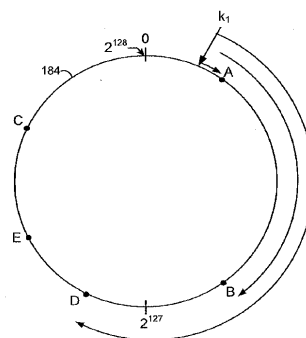
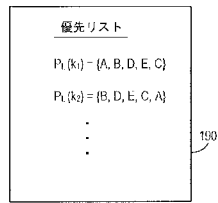


FIG. 10

【図 1 1】



【図 1 2】

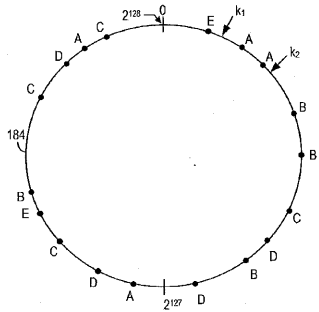
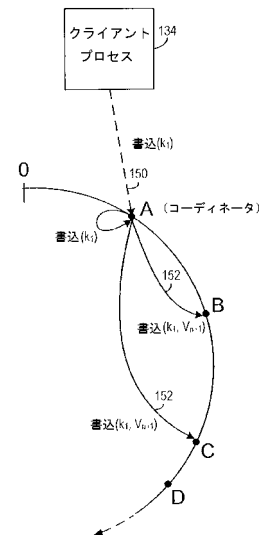
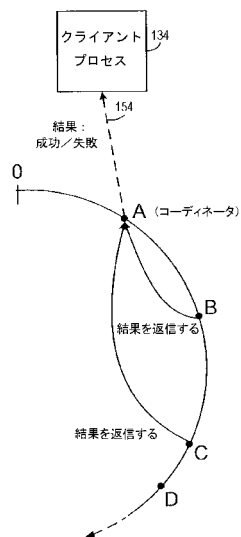


FIG. 12

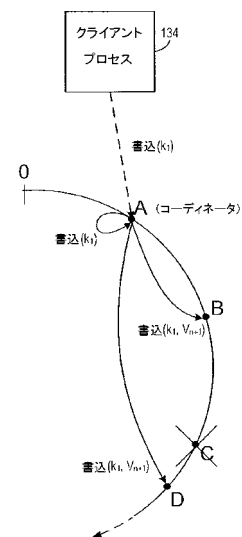
【図 1 3 A】



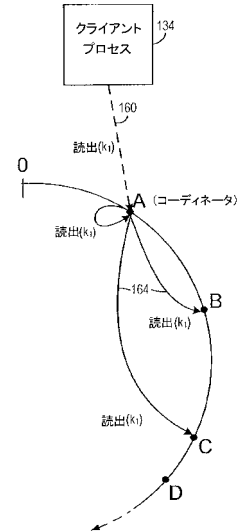
【図 1 3 B】



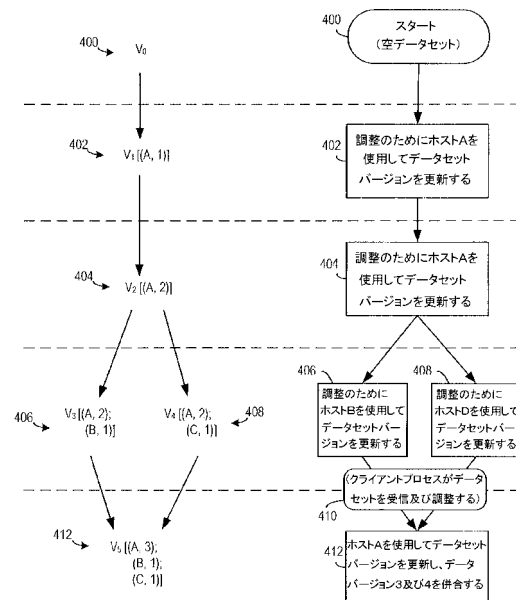
【図 1 4 A】



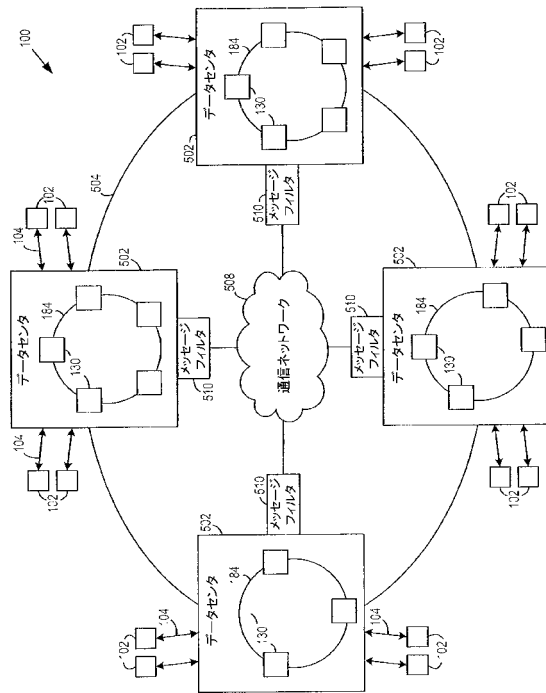
【 図 1 5 A 】



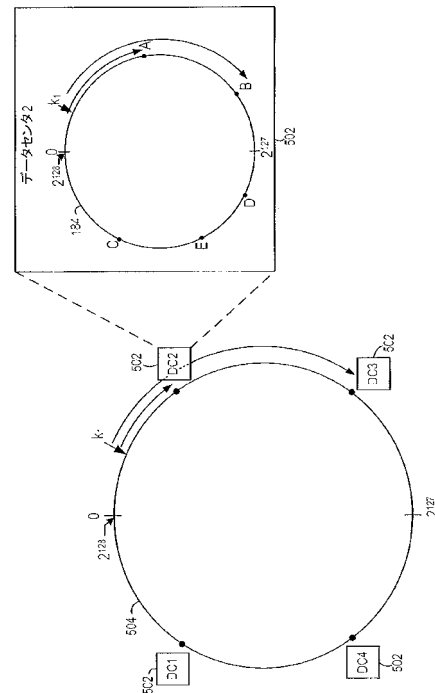
【 図 1 6 】



【図 17】



【図 18】



【図 19】

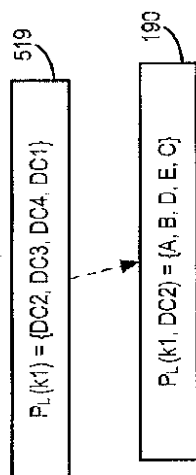
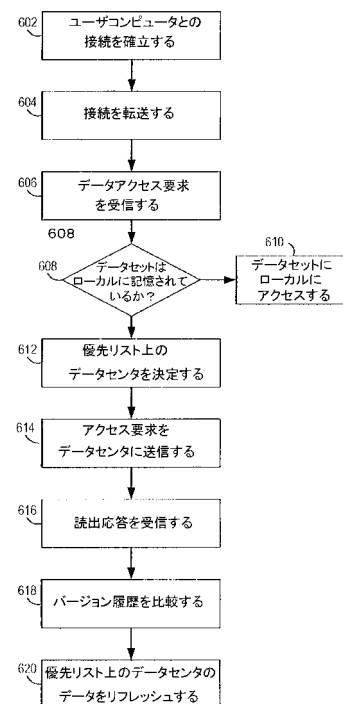
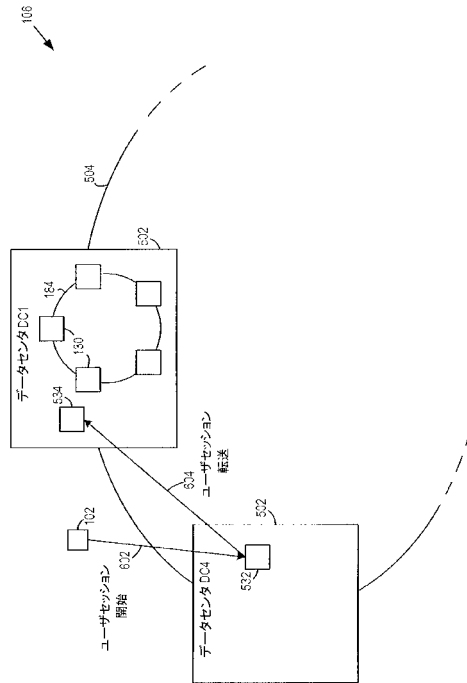


FIG. 19

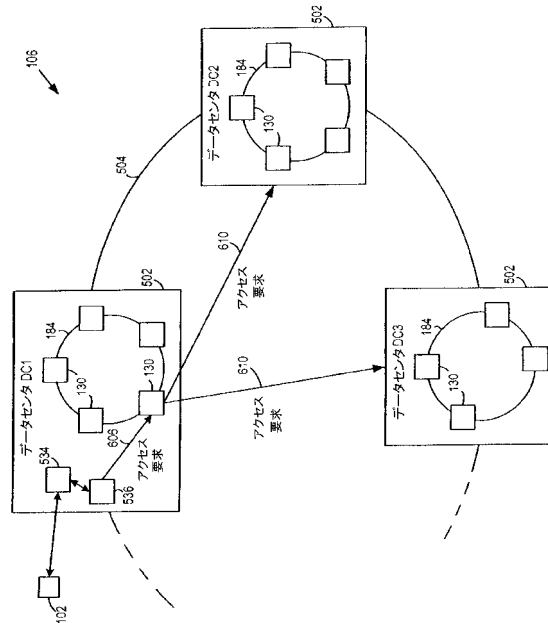
【図 20】



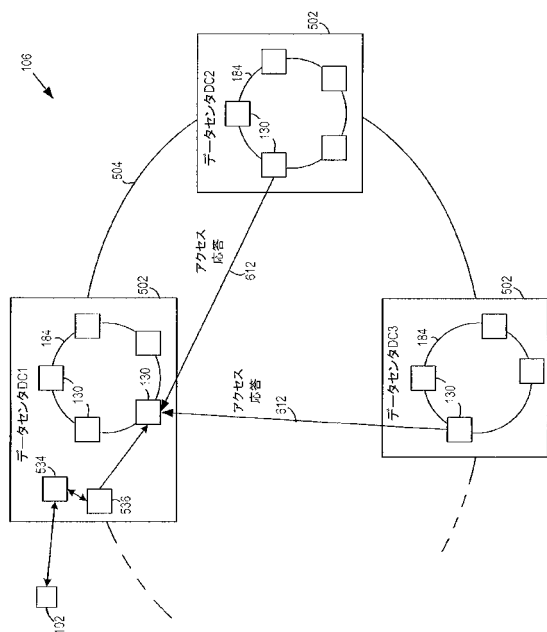
【 図 2 1 】



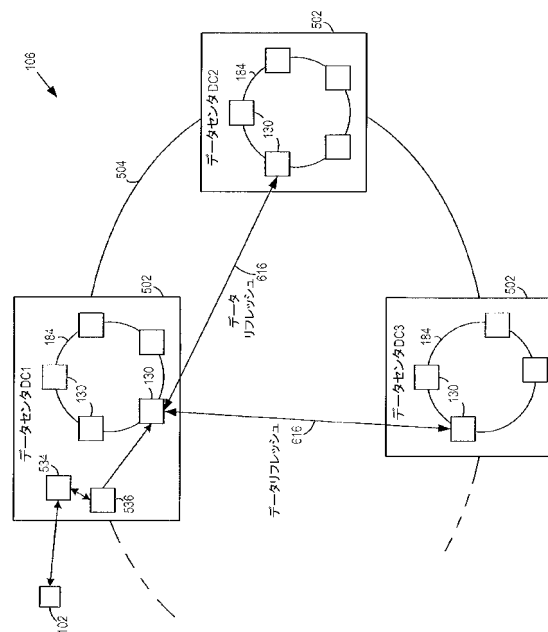
【 図 2 2 】



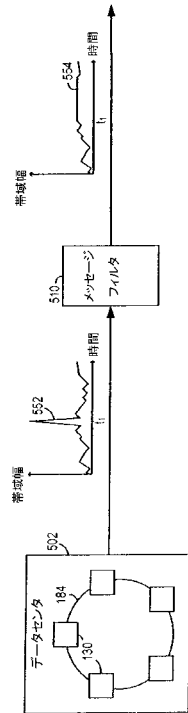
【圖 23】



【 図 2 4 】



【図 25】



フロントページの続き

- (74)代理人 100135633
弁理士 二宮 浩康
- (74)代理人 100114890
弁理士 アインゼル・フェリックス＝ラインハルト
- (72)発明者 ピーター ヴォスホール
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0
- (72)発明者 スワミナサン シヴァスブラマニアン
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0
- (72)発明者 ジョゼッペ デカンディア
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0
- (72)発明者 デニス ハストルン
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0
- (72)発明者 アヴィナシュ ラクシュマン
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0
- (72)発明者 アレックス ピルチン
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0
- (72)発明者 イヴァン ディー ロセロ
アメリカ合衆国 ワシントン シアトル トウエルヴス アヴェニュー サウス 1 2 0 0 -
スイート 1 2 0 0

審査官 原 秀人

- (56)参考文献 特開 2 0 0 6 - 1 9 5 7 4 4 (J P , A)
特開 2 0 0 4 - 1 5 7 7 9 3 (J P , A)
米国特許出願公開第 2 0 0 5 / 0 1 2 0 0 2 5 (U S , A 1)
国際公開第 2 0 0 6 / 0 5 3 5 8 0 (W O , A 1)
国際公開第 2 0 0 5 / 0 4 3 3 2 3 (W O , A 1)
米国特許第 0 5 9 2 4 0 9 6 (U S , A)
米国特許第 0 6 5 9 7 7 0 0 (U S , B 1)
DABEK F ET AL. , Wide-area cooperative storage with CFS , 2 0 0 1 年 1 0 月 2 1 日 , U R L
 , http://ftp.lcdf.org/papers/cfs:sosp01/cfs_sosp.pdf

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 2 / 0 0
G 0 6 F 1 3 / 0 0