

(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2016-510458
(P2016-510458A)

(43) 公表日 平成28年4月7日(2016.4.7)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 304F	5B005
G06F 12/10 (2016.01)	G06F 3/06 301Z	
	G06F 3/06 301F	
	G06F 12/10 501Z	

審査請求 未請求 予備審査請求 未請求 (全 21 頁)

(21) 出願番号 特願2015-552744 (P2015-552744)
 (86) (22) 出願日 平成26年1月8日 (2014.1.8)
 (85) 翻訳文提出日 平成27年8月18日 (2015.8.18)
 (86) 国際出願番号 PCT/US2014/010690
 (87) 国際公開番号 W02014/110137
 (87) 国際公開日 平成26年7月17日 (2014.7.17)
 (31) 優先権主張番号 61/751, 142
 (32) 優先日 平成25年1月10日 (2013.1.10)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 14/046, 872
 (32) 優先日 平成25年10月4日 (2013.10.4)
 (33) 優先権主張国 米国 (US)

(71) 出願人 513076589
 ピュア・ストレージ・インコーポレイテッド
 アメリカ合衆国・94041・カリフォルニア州・マウンテンビュー・カストロストリート・650・スイート・220
 (74) 代理人 100064621
 弁理士 山川 政樹
 (74) 代理人 100098394
 弁理士 山川 茂樹
 (72) 発明者 コルグローヴ, ジョン
 アメリカ合衆国・94024・カリフォルニア州・ロスアルトス・ヴィスタグラウンデアヴェニュー・722

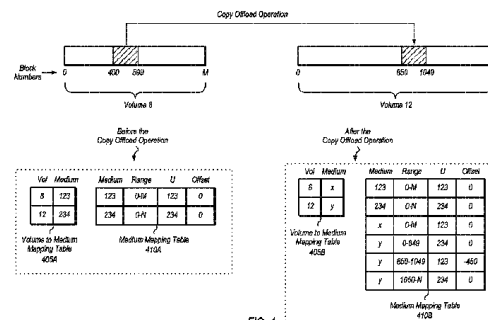
最終頁に続く

(54) 【発明の名称】 ストレージシステム内でのコピーの実行

(57) 【要約】

コピーオフロード動作を実行するためのシステム及び方法。(第1の媒体を指す)第1のボリュームから(第2の媒体を指す)第2のボリューム)へのコピーオフロード動作を要求されると、このコピーオフロード動作は、コピーされるデータにアクセスすることなく実行される。第3の媒体が生成され、第1の媒体はこの第3の媒体の下層媒体として記録される。第1のボリュームは第3の媒体を再び指す。また第4の媒体が生成され、第2のボリュームは第4の媒体を再び指し、第2の媒体は第4の媒体の標的範囲の下層媒体として記録される。第4の媒体の他の全ての範囲はその下層媒体として第2の媒体を有する。

【選択図】 図4



【特許請求の範囲】

【請求項 1】

1つ又は複数のストレージデバイス；及び

前記1つ又は複数のストレージデバイスに接続されたストレージコントローラを備える、コンピュータシステムであって、

前記ストレージコントローラは、第1の媒体に関連する第1のボリュームからの第1の範囲のデータを、第2の媒体に関連する第2のボリュームの第2の範囲へコピーするための要求の受信に応答して；

第3の媒体を生成し、前記第1の媒体が前記第3の媒体の下層にあることのインジケータを記憶し、前記第1のボリュームを前記第3の媒体に関連付けるよう；

第4の媒体を生成し、前記第2の媒体が前記第4の媒体の下層にあることのインジケータを記憶し、前記第2のボリュームを前記第4の媒体に関連付けるよう；及び

前記第1の媒体内の前記第1の範囲のデータが、前記第4の媒体の前記第2の範囲の下層にあることのインジケータを記憶するよう構成される、コンピュータシステム。

【請求項 2】

前記第1の範囲のデータをコピーするための前記要求は、前記第1の範囲のデータにアクセスする必要なしに、前記ストレージコントローラによって満たされる、請求項1に記載のコンピュータシステム。

【請求項 3】

前記ストレージコントローラは更に、前記第3の媒体の生成を、前記第1のボリュームの前記第1の範囲を標的とする書き込み要求を受信するまで遅延させるよう構成される、請求項1に記載のコンピュータシステム。

【請求項 4】

前記ストレージコントローラは更に、前記第4の媒体の生成を、前記第2のボリュームの前記第2の範囲を標的とする書き込み要求を受信するまで遅延させるよう構成される、請求項3に記載のコンピュータシステム。

【請求項 5】

前記第1の範囲は、前記第1のボリューム内に第1のオフセットで配置され、

前記第2の範囲は、前記第2のボリューム内に第2のオフセットで配置され、

前記ストレージコントローラは、前記第2のオフセットと前記第1のオフセットとの差を、媒体マッピングテーブル内のエントリに記憶するよう構成され、

前記エントリは前記第4の媒体の前記第2の範囲に対応する、請求項1に記載のコンピュータシステム。

【請求項 6】

前記ストレージコントローラは更に；

前記第1の媒体が前記第3の媒体の下層にあることのインジケータを記憶するのに対応して、前記第1の媒体が読み出し専用であることのインジケータを記憶するよう；及び

前記第2の媒体が前記第4の媒体の下層にあることのインジケータを記憶するのに対応して、前記第2の媒体が読み出し専用であることのインジケータを記憶するよう構成される、請求項1に記載のコンピュータシステム。

【請求項 7】

前記第1の範囲のデータは第5の媒体内に記憶され、

前記第5の媒体は前記第1の媒体の下層にある、請求項1に記載のコンピュータシステム。

【請求項 8】

第1のボリュームからの第1の範囲のデータを、第2のボリュームの第2の範囲へコピーするための要求の受信に応答して；

第3の媒体を生成するステップ；

第1の媒体が前記第3の媒体の下層にあることのインジケータを記憶し、ここで前記第

10

20

30

40

50

- 1 のボリュームの前記第 1 の範囲は前記第 1 の媒体に関連する、ステップ；
 前記第 1 のボリュームを前記第 3 の媒体と関連付けるステップ；
 第 4 の媒体を生成するステップ；
 前記第 1 の媒体内の前記第 1 の範囲のデータが前記第 4 の媒体の下層にあることのインジケータを記憶するステップ；
 第 2 の媒体が、前記第 4 の媒体の他の全てのデータの下層にあることの 1 つ又は複数のインジケータを記憶し、ここで前記第 2 のボリュームは前記第 2 の媒体に関連する、ステップ；及び
 前記第 2 のボリュームを前記第 4 の媒体と関連付けるステップ
 を含む、ストレージシステムにおいて使用するための方法。 10
- 【請求項 9】
 前記第 1 の範囲のデータをコピーするための前記要求を、前記第 1 の範囲のデータにアクセスすることなく満たすステップを更に含む、請求項 8 に記載の方法。
- 【請求項 10】
 前記第 3 の媒体の生成を、前記第 1 のボリュームの前記第 1 の範囲を標的とする書き込み要求を受信するまで遅延させるステップを更に含む、請求項 8 に記載の方法。
- 【請求項 11】
 前記第 4 の媒体の生成を、前記第 2 のボリュームの前記第 2 の範囲を標的とする書き込み要求を受信するまで遅延させるステップを更に含む、請求項 10 に記載の方法。 20
- 【請求項 12】
 前記第 1 の範囲は、前記第 1 のボリューム内に第 1 のオフセットで配置され、
 前記第 2 の範囲は、前記第 2 のボリューム内に第 2 のオフセットで配置され、
 前記方法は、前記第 2 のオフセットと前記第 1 のオフセットとの差を、媒体マッピングテーブル内のエントリに記憶するステップを更に含む
 前記エントリは前記第 4 の媒体の前記第 2 の範囲に対応する、請求項 8 に記載の方法。
- 【請求項 13】
 前記第 1 の媒体が前記第 3 の媒体の下層にあることのインジケータを記憶するのに対応して、前記第 1 の媒体が読み出し専用であることのインジケータを記憶するステップ；及び
 前記第 2 の媒体が前記第 4 の媒体の下層にあることのインジケータを記憶するのに対応して、前記第 2 の媒体が読み出し専用であることのインジケータを記憶するステップ
 を更に含む、請求項 8 に記載の方法。 30
- 【請求項 14】
 前記第 1 の範囲のデータは第 5 の媒体内に記憶され、
 前記第 5 の媒体は前記第 1 の媒体の下層にある、請求項 13 に記載の方法。
- 【請求項 15】
 プログラム命令を記憶する非一時的コンピュータ可読ストレージ媒体であって、
 第 1 のボリュームからの第 1 の範囲のデータを、第 2 のボリュームの第 2 の範囲へコピーするための要求の受信に対応して：
 第 3 の媒体を生成し； 40
 第 1 の媒体が前記第 3 の媒体の下層にあることのインジケータを記憶し、ここで前記第 1 のボリュームは前記第 1 の媒体に関連し；
 前記第 1 のボリュームを前記第 3 の媒体と関連付け；
 第 4 の媒体を生成し；
 前記第 1 の媒体内の前記第 1 の範囲のデータが前記第 4 の媒体の下層にあることのインジケータを記憶し；
 第 2 の媒体が、前記第 4 の媒体の他の全てのデータの下層にあることの 1 つ又は複数のインジケータを記憶し、ここで前記第 2 のボリュームは前記第 2 の媒体に関連し；及び
 前記第 2 のボリュームを前記第 3 の媒体と関連付ける
 ために、前記プログラム命令をプロセッサによって実行できる、コンピュータ可読ストレ 50

ージ媒体。

【請求項 16】

前記第 1 の範囲のデータをコピーするための前記要求を、前記第 1 の範囲のデータにアクセスすることなく満たすために、前記プログラム命令をプロセッサによって更に実行できる、請求項 15 に記載のコンピュータ可読ストレージ媒体。

【請求項 17】

前記第 3 の媒体の生成を、前記第 1 のボリュームの前記第 1 の範囲を標的とする書き込み要求を受信するまで遅延させるために、前記プログラム命令をプロセッサによって更に実行できる、請求項 15 に記載のコンピュータ可読ストレージ媒体。

【請求項 18】

前記第 4 の媒体の生成を、前記第 2 のボリュームの前記第 2 の範囲を標的とするデータ要求を受信するまで遅延させるために、前記プログラム命令をプロセッサによって更に実行できる、請求項 17 に記載のコンピュータ可読ストレージ媒体。

【請求項 19】

前記第 1 の範囲は、前記第 1 のボリューム内に第 1 のオフセットで配置され、
前記第 2 の範囲は、前記第 2 のボリューム内に第 2 のオフセットで配置され、
前記第 2 のオフセットと前記第 1 のオフセットとの差を、媒体マッピングテーブル内のエンTRIESに記憶するために、前記プログラム命令をプロセッサによって更に実行でき、
前記エンTRIESは前記第 4 の媒体の前記第 2 の範囲に対応する、請求項 15 に記載のコンピュータ可読ストレージ媒体。

【請求項 20】

前記第 1 の媒体が前記第 3 の媒体の下層にあることのインジケータを記憶するのに応答して、前記第 1 の媒体が読み出し専用であることのインジケータを記憶し；及び
前記第 2 の媒体が前記第 4 の媒体の下層にあることのインジケータを記憶するのに応答して、前記第 2 の媒体が読み出し専用であることのインジケータを記憶するために、前記プログラム命令をプロセッサによって更に実行できる、請求項 15 に記載のコンピュータ可読ストレージ媒体。

【請求項 21】

1 つ又は複数のストレージデバイス；及び
前記 1 つ又は複数のストレージデバイスに接続されたストレージコントローラ
を備える、コンピュータシステムであって、
前記ストレージコントローラは、第 1 のボリューム内の第 1 の既存の媒体に関連する第 1 の位置からの第 1 のデータを、第 2 の位置へコピーするための要求の受信に応答して：
第 1 の新規の媒体を生成し；
前記第 1 の既存の媒体が前記第 1 の新規の媒体の下層にあることの第 1 のインジケータを記憶し；
前記第 1 のボリュームを前記第 1 の新規の媒体と関連付け；
前記第 2 の位置が前記第 1 のボリューム内にある場合、前記第 1 の媒体内の前記第 1 の位置が前記第 1 の新規の媒体内の前記第 2 の位置の下層にあることの第 2 のインジケータを記憶し；
前記第 2 の位置が、第 2 の既存の媒体に関連する第 2 のボリューム内にある場合：
第 2 の新規の媒体を生成し、
前記第 2 の既存の媒体が前記第 2 の新規の媒体の下層にあることの第 3 のインジケータを記憶し；
前記第 1 のデータが前記第 2 の新規の媒体の下層にあることの第 4 のインジケータを記憶し；
前記第 2 のボリュームを前記第 2 の新規の媒体と関連付ける
よう構成される、コンピュータシステム。

【請求項 22】

1 つ又は複数のストレージデバイス；及び

10

20

30

40

50

前記 1 つ又は複数のストレージデバイスに接続されたストレージコントローラを備える、コンピュータシステムであって、

前記ストレージコントローラは、第 1 のボリューム内の第 1 の既存の媒体に関連する第 1 の位置からの第 1 のデータを、第 2 の位置へコピーするための要求の受信に応答して：
第 1 の新規の媒体を生成し；

前記第 1 の既存の媒体が前記第 1 の新規の媒体の下層にあることの第 1 のインジケータを記憶し；

前記第 1 のボリュームを前記第 1 の新規の媒体と関連付け；

前記第 1 の媒体内の前記第 1 の位置が前記第 1 の新規の媒体内の前記第 2 の位置の下層にあることの第 2 のインジケータを記憶する

よう構成される、コンピュータシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステム内でのコピーオフロード動作の実行に関する。

【背景技術】

【0002】

データ保守は、製品開発、人的資源、財務又はその他のいずれに関する業務であっても、その操作の重要な部分である。このようなデータ保守の一部として、多様な複数の理由のうちのいずれかにより、データのコピーを生成する必要がある場合がある。コピーを生成するための 1 つのアプローチは、（例えばネットワークを介してストレージシステムに接続された）クライアントがコピープロセスにアクティブに従事することを伴う。このアプローチでは、1 つ又は複数のクライアントプロセッサが、コピープロセス全体を通してストレージシステムに搬送されるトランザクションを開始させ、管理する。従ってクライアントプロセッサはコピー動作によってビジー状態となり、一般に他の作業に従事できなくなる。データをコピーするための別のアプローチは、コピープロセスの管理を別の 1 つ又は複数のプロセッサにオフロードすることを伴う。例えばコピープロセスをストレージシステムにオフロードし、このストレージシステムがコピープロセスを管理してよい。このようなアプローチは「コピーオフロード」と呼ばれることがある。この方法では、クライアントは一般に、コピーを実行している間に他の作業を実行できるよう自由である。コピーオフロード動作は、ストレージシステムがデータ（場合によっては大量のデータ）を、単一のボリューム内又は 2 つのボリューム間において、ある位置から別の位置へとコピーする動作を意味する。

【0003】

以上に鑑みて、コピーオフロード動作を効率的に実行するためのシステム及び方法に対する需要が存在する。

【発明の概要】

【発明が解決しようとする課題】

【0004】

コピーオフロード動作を実行するためのシステム及び方法の様々な実施形態を考える。

【課題を解決するための手段】

【0005】

ストレージシステムは、ストレージコントローラ及び 1 つ又は複数のストレージデバイスを含んでよい。このストレージシステムは、1 つ又は複数のホストクライアントシステムに接続してよい。一実施形態では、ストレージコントローラはボリューム及び媒体を利用して、ストレージシステム内に記憶されているクライアントデータを追跡する。媒体はデータの論理グループ化として定義され、各媒体は、データの論理グループ化を識別するための識別子を有する。ストレージコントローラは、ボリューム - 媒体マッピングテーブルを保守することによって各ボリュームを単一の媒体にマッピングでき、この媒体をそのボリュームのアンカー媒体と呼ぶ。各媒体は、いずれの個数の他の媒体に対してマッピン

10

20

30

40

50

ができ、ストレージコントローラはまた、媒体マッピングテーブルを保守することによってマッピングを媒体間で追跡できる。

【0006】

ホストシステムが、データを第1のボリュームから第2のボリュームへ（又は第1のボリューム内の第1の位置から第1のボリューム内の第2の位置へ）データをコピーするためのコピーオフロード動作を要求すると、ストレージコントローラは、ボリューム-媒体マッピングテーブルを操作するだけでコピーオフロード動作を実行でき、この媒体マッピングテーブルはコピー中のデータにアクセスすることはない。その結果、コピーオフロード動作を迅速に、かつストレージシステムのリソースの利用を最低限とした状態で実行できる。

10

【0007】

これらの及びその他の実施形態は、以下の説明及び添付の図面を考察することによって明らかになるであろう。

【図面の簡単な説明】

【0008】

【図1】図1は、ストレージシステムの一実施形態を示す概略ブロック図である。

【図2】図2は、媒体の有向非巡回グラフ(directed acyclic graph: DAG)の一実施形態の概略ブロック図である。

【図3】図3は、媒体マッピングテーブルの一実施形態を示す。

【図4】図4は、コピーオフロード動作の一実施形態を示す。

20

【図5】図5は、コピーオフロード動作の別の実施形態の概略ブロック図である。

【図6】図6は、コピーオフロード動作を実行するための方法の一実施形態を示す概略フロー図である。

【図7】図7は、コピーオフロード動作を実行するための別の方法の一実施形態を示す概略フロー図である。

【発明を実施するための形態】

【0009】

本発明は様々な修正例及び代替形態を許容するものであるが、例として複数の具体的実施形態を図面に示し、また本明細書において詳細に説明する。しかしながら、図面及びこれに対する詳細な説明は、本出願において開示されている特定の形態に本発明を限定することを意図したものではなく、その反対に、本発明は、添付の請求項によって定義される本発明の精神及び範囲内のあらゆる修正例、均等物、代替例を包含できることを理解されたい。

30

【0010】

以下の説明では、本発明の完全な理解を提供するために、多数の具体的な詳細を挙げる。しかしながら、これらの具体的な詳細を用いずに本発明を実施してもよいことを、当業者は認識するべきである。例えば公知の回路、構造、信号、コンピュータプログラム命令、技術については、本発明を不明瞭にするのを回避するために、詳細には示していない。

【0011】

ここで図1を参照すると、ストレージシステム100の一実施形態の概略ブロック図が示されている。ストレージシステム100は、ストレージコントローラ110及びストレージデバイス群130、140を含んでよく、上記ストレージデバイス群130、140はいずれの個数のストレージデバイス群（又はデータストレージアレイ）の代表例である。図示したように、ストレージデバイス群130はストレージデバイス135A~Nを含み、これらはいずれの個数及びタイプのストレージデバイス（例えばソリッドステートドライブ(solid-state drive: SSD)）の代表例である。ストレージコントローラ110はクライアントコンピュータシステム125に直接接続してよく、またストレージコントローラ110は、ネットワーク120を介して、クライアントコンピュータシステム115から離間した状態でクライアントコンピュータシステム115に接続してよい。クライアント115、125は、システム100においてデータを記憶しデ

40

50

ータにアクセスするためにストレージコントローラ 110 を利用できるいずれの個数のクライアントの代表例である。なお、いくつかのシステムは、ストレージコントローラ 110 に直接又は離間した状態で接続された単一のクライアントしか含まなくてもよい。

【0012】

ストレージコントローラ 110 は、ストレージデバイス 135 A ~ N へのアクセスを提供するよう構成されたソフトウェア及び/又はハードウェアを含んでよい。ストレージデバイス群 130、140 から離間したものとしてストレージコントローラ 110 を図示しているが、いくつかの実施形態では、ストレージコントローラ 110 をストレージデバイス群 130、140 のうちの一方又はそれぞれの中に配置してよい。ストレージコントローラ 110 は、ベースオペレーティングシステム (operating system: OS)、ボリュームマネージャ、本明細書で開示される様々な技術を実装するための追加の制御論理を含んでよく、又はこれらに接続してよい。

10

【0013】

ストレージコントローラ 110 は実施形態に応じて、いずれの個数のプロセッサを含んでよく、及び/又はいずれの個数のプロセッサ上で実行されてよく、また、単一のホストコンピューティングデバイスを含んでよく、及び/若しくは単一のホストコンピューティングデバイス上で実行されてよく、又は複数のホストコンピューティングデバイスに亘って分散されていてよい。いくつかの実施形態では、ストレージコントローラ 110 は一般に、1つ若しくは複数のファイルサーバ及び/若しくはブロックサーバを含むか、又は1つ若しくは複数のファイルサーバ及び/若しくはブロックサーバ上で実行されてよい。ストレージコントローラ 110 は、デバイスの不良又はデバイス内のストレージ位置の不良によるデータの損失を防止するために、デバイス 135 A ~ N に亘ってデータを複製するための様々な技術のうちのいずれを使用してよい。ストレージコントローラ 110 は、共通のデータを重複排除することによって、デバイス 135 A ~ N に記憶されるデータ量を削減するための、様々な重複排除、圧縮又はその他の技術のうちのいずれを利用してよい。

20

【0014】

ストレージコントローラ 110 はまた、システム 100 内でスナップショットを生成してこれを管理するよう構成してよい。媒体のセットは、ストレージコントローラ 110 によって記録され、保守される。媒体のほとんどは、特定のボリュームが使用中の最新の媒体等の1つ又は複数の選択された媒体を除いて、読み出し専用であってよい。各媒体は論理的に、媒体内の全てのブロックを含む。しかしながら、媒体が生成された時点から媒体が閉鎖された時点までに変化したブロックのみが保存され、これらのブロックに対するマッピングを上記媒体を用いて保守してもよい。

30

【0015】

様々な実施形態では、ストレージコントローラ 110 が複数のマッピングテーブルを保守してよい。これらのマッピングテーブルは、媒体マッピングテーブル及びボリューム - 媒体間マッピングテーブルを含んでよい。これらのテーブルを利用して、媒体と下層媒体との間のマッピング及びボリュームと媒体との間のマッピングを記録及び保守してよい。ストレージコントローラ 110 はまた、複数のエントリを有するアドレス翻訳テーブルを含んでよく、各エントリは、対応するデータ成分に関する仮想 - 物理間マッピングを保持する。このマッピングテーブルを用いて、クライアントコンピュータシステム 115、125 からストレージデバイス 135 A ~ N 内の物理的位置へ、論理読み書き要求をマッピングできる。受信した読み書き要求に対応する探索動作の間に、所定の媒体に関連するマッピングから「物理」ポインタ値を読み出してよい。用語「マッピング (mapping)」は、所定の媒体 ID 及びブロック番号を物理ポインタ値に変換するアドレス翻訳マッピングテーブルの1つ又は複数のエントリとして定義される。続いてこの物理ポインタ値を用いて、ストレージデバイス 135 A ~ N 内の物理的位置を配置してよい。なお、物理ポインタ値を用いて、ストレージデバイス 135 A ~ N のうちの所定のストレージデバイス内の別のマッピングテーブルにアクセスしてもよい。その結果、物理ポインタ値と標的

40

50

ストレージ位置との間に1つ又は複数のレベル間接参照を存在させることができる。

【0016】

なお代替実施形態では、クライアントコンピュータ、ストレージコントローラ、ネットワーク、ストレージデバイス群、データストレージデバイスの個数及びタイプは、図1に示したものに限定されない。様々な時点において、1つ又は複数のクライアントがオフラインで動作できる。更に動作中、ユーザがシステム100に対して接続、接続解除、再接続を行うのに従って、個々のクライアントコンピュータ接続タイプを変更してよい。更に本明細書に記載のシステム及び方法は、直接接続型ストレージシステム又はネットワーク接続型ストレージシステムに適用してよく、また本明細書に記載の方法の1つ又は複数の態様を実行するよう構成されたホストオペレーティングシステムを含んでよい。多数のこのような代替例が可能であり、考察の対象となる。

10

【0017】

ネットワーク120は：無線接続；直接ローカルエリアネットワーク(local area network：LAN)接続；インターネット、ルータ、ストレージエリアネットワーク、イーサネット(登録商標)等の広域ネットワーク(wide area network：WAN)接続を含む多様な技術を利用してよい。ネットワーク120は、無線式であってもよい1つ又は複数のLANを備えてよい。ネットワーク120は更に、リモートダイレクトメモリアクセス(remote direct memory access：RDMA)ハードウェア及び/若しくはソフトウェア、伝送制御プロトコル/インターネット・プロトコル(transmission control protocol/internet protocol：TCP/IP)ハードウェア及び/若しくはソフトウェア、ルータ、リピータ、スイッチ並びに/又はグリッド並びに/又はその他のものを含んでよい。ファイバチャネル、ファイバチャネルオーバーイーサネット(Fibre Channel over Ethernet：FCoE)、iSCSI等のプロトコルをネットワーク120において使用してよい。ネットワーク120は、伝送制御プロトコル(Transmission Control Protocol：TCP)及びインターネット・プロトコル(Internet Protocol：IP)、即ちTCP/IP等の、インターネット用に使用される一連の通信プロトコルとインタフェース接続してよい。

20

【0018】

クライアントコンピュータ115、125は、デスクトップパーソナルコンピュータ(personal computer：PC)、サーバ、サーバファーム、ワークステーション、ラップトップ、ハンドヘルドコンピュータ、サーバ、パーソナルデジタルアシスタント(personal digital assistant：PDA)、スマートフォン等の、いずれの個数の据置型又は移動体コンピュータの代表例である。一般にクライアントコンピュータシステム115、125は、1つ又は複数のプロセッサコアを備える1つ又は複数のプロセッサを含む。各プロセッサコアは、事前に定義された汎用命令セットに従って命令を実行するための回路構成を含む。例えばx86命令セットアーキテクチャを選択してよい。あるいはARM(登録商標)、Alpha(登録商標)、PowerPC(登録商標)、SPARC(登録商標)又は他のいずれの汎用命令セットアーキテクチャを選択してよい。プロセッサコアは、データ及びコンピュータプログラム命令のためのキャッシュメモリサブシステムにアクセスしてよい。上記キャッシュサブシステムは、ランダムアクセスメモリ(random access memory：RAM)及びストレージデバイスを備えるメモリ階層に接続してよい。

30

40

【0019】

ここで図2を参照すると、媒体の有向非巡回グラフ(DAG)200を示すブロック図が示されている。また図2は、各ボリュームに関して、ストレージシステムによって使用されている間にボリュームがどの媒体にマッピングされるかを示す、ボリューム-媒体間マッピングテーブル205も示す。ボリュームはグラフ200への考えられるポイントであり得る。

50

【0020】

本明細書で使用される用語「媒体 (medium)」は、データの論理グループ化として定義される。媒体は、データの論理グループ化を識別するための対応する識別子を有してよい。各媒体はまた、コンテンツ位置、重複排除エントリ及びその他の情報に対する論理ブロック番号のマッピングを含むか、又はこれに関連してよい。一実施形態では、媒体の識別子をストレージコントローラが使用してよいが、媒体の識別子はユーザ管理下でなくてよい。ユーザ (又はクライアント) は、ボリュームIDに付随するデータ要求を送信して、この要求がどのデータを標的とするかを特定してよく、ストレージコントローラはボリュームIDを媒体IDに対してマッピングして、上記要求を処理する際に媒体IDを使用してよい。

10

【0021】

用語「媒体」を用語「ストレージ媒体 (storage medium)」又は「コンピュータ可読ストレージ媒体 (computer readable storage medium)」と混同してはならない。ストレージ媒体は、データを記憶するために利用される実際の物理デバイス (例えばSSD、HDD) として定義される。コンピュータ可読ストレージ媒体 (又は非一時的コンピュータ可読ストレージ媒体) は、プロセッサ又はその他のハードウェアデバイスが実行できるプログラム命令を記憶するよう構成された物理ストレージ媒体として定義される。本明細書に記載の方法及び/又はメカニズムを実装する様々なタイプのプログラム命令を、コンピュータ可読媒体上で搬送又は記憶してよい。プログラム命令を記憶するよう構成された多数のタイプの媒体利用手可能であり、これにはハードディスク、フロッピー (登録商標) ディスク、CD-ROM、DVD、フラッシュメモリ、プログラマブルROM (Programmable ROM: PROM)、ランダムアクセスメモリ (RAM)、他の様々な形態の揮発性又は不揮発性ストレージが含まれる。

20

【0022】

また、用語「ボリューム - 媒体間マッピングテーブル (volume to medium mapping table)」は、単なる単一のテーブルではなく複数のテーブルを指してよいことにも留意されたい。同様に用語「媒体マッピングテーブル (medium mapping table)」も、単なる単一のテーブルではなく複数のテーブルを指してよい。更に、ボリューム - 媒体間マッピングテーブル205は、ボリューム - 媒体間マッピングテーブルの一例にすぎないことにも留意されたい。他のボリューム - 媒体間マッピングテーブルは、異なる個数のボリュームのための異なる個数のエントリを有してよい。

30

【0023】

各媒体を、3つの結合されたボックスとしてグラフ200に示し、左のボックスは媒体IDを示し、中央のボックスは下層媒体を示し、右のボックスは媒体のステータス (RO: 読み出し専用 (read-only)) 又は (RW: 読み書き (read-write)) を示す。グラフ200内では、媒体はその下層媒体を指す。例えば媒体20は媒体12を指し、これによって媒体12が媒体20の下層媒体であることを表す。また媒体12は媒体10を指し、媒体10は媒体5を指し、媒体5は媒体1を指す。いくつかの媒体は、2つ以上の高次の媒体に対する下層媒体となる。例えば3つの別個の媒体 (12、17、11) は媒体10を指し、2つの別個の媒体 (18、10) は媒体5を指し、2つの別個の媒体 (6、5) は媒体1を指す。少なくとも1つの高次媒体に対する下層媒体となる媒体はそれぞれ、読み出し専用のステータスを有する。

40

【0024】

グラフ200の左下の媒体のセットは、線形セットの例である。グラフ200に示すように、まず媒体3が生成され、続いてスナップショットが実行され、その結果媒体3は安定する (即ちこの時点以降、媒体3内の所定のブロックの探索の結果が常に同一の値を返すことになる)。媒体3を下層媒体として用いて、媒体7が生成される。媒体3が安定した後に書かれるいずれのブロックは、媒体7内にあるものとして標識される。媒体7に対

50

する探索は、媒体 7 内に値が見つかる場合は媒体 7 から値を返すが、媒体 7 内にブロックが見つからない場合は媒体 3 を探索することになる。その後媒体 7 のスナップショットを実行し、媒体 7 は安定となり、媒体 1 4 が生成される。媒体 1 4 内のブロックの探索により、媒体 7、続いて媒体 3 を検査して、標的論理ブロックを見つける。最後に媒体 1 4 のスナップショットを実行し、媒体 1 4 は安定となり、媒体 1 5 が生成される。グラフ 2 0 0 のこの時点において媒体 1 4 は、媒体 1 5 へと向かうボリューム 1 0 2 への書き込みによって安定となる。

【 0 0 2 5 】

ボリューム - 媒体間マッピングテーブル 2 0 5 は、ユーザ管理下のボリュームを媒体に対してマッピングする。各ボリュームは、アンカー媒体としても知られる単一の媒体に対してマッピングしてよい。このアンカー媒体は、他の全ての媒体と同様に、それ固有の探索を処理してよい。複数のボリュームが依存する媒体（例えば媒体 1 0）は、それ自体の複数のブロックを、これらブロックが依存するボリュームとは別個に追跡する。各媒体をブロックの範囲に分解してもよく、各範囲を媒体の D A G 2 0 0 内で別個に処理してよい。

10

【 0 0 2 6 】

ここで図 3 を参照すると、媒体マッピングテーブル 3 0 0 の一実施形態を示す。媒体マッピングテーブル 3 0 0 のいずれの部分又は全体は、ストレージコントローラ 1 1 0 及び/又はストレージデバイス 1 3 5 A ~ N のうちの 1 つ若しくは複数に記憶してよい。ボリューム識別子 (I D) を使用してボリューム - 媒体間マッピングテーブル 2 0 5 にアクセスし、このボリューム I D に対応する媒体 I D を決定してよい。次にこの媒体 I D を使用して媒体マッピングテーブル 3 0 0 にアクセスしてよい。なおテーブル 3 0 0 は媒体マッピングテーブルの単なる一例であり、他の実施形態では、異なる個数のエントリを有する異なる媒体マッピングテーブルを利用してよい。更に他の実施形態では、媒体マッピングテーブルは他の属性を含んでよく、図 3 に示したものと異なる方法で整理されてよい。

20

【 0 0 2 7 】

テーブル 3 0 0 の左列に示すように、媒体 I D によって各媒体を識別してよい。テーブル 3 0 0 の各エントリに範囲属性も含んでよく、この範囲はデータブロックに関するものであってよい。データのブロックのサイズ（例えば 4 K B、8 K B）は、実施形態に応じて変化し得る。媒体は複数の範囲に分解してよく、媒体の各範囲を、固有の属性及びマッピングを有する独立した媒体であるかのように処理してよい。例えば媒体 I D 2 は 2 つの別個の範囲を有する。媒体 I D 2 の範囲 0 ~ 9 9 は、テーブル 3 0 0 内で媒体 I D 2 の範囲 1 0 0 ~ 9 9 9 のためのエントリから離間したエントリを有する。

30

【 0 0 2 8 】

媒体 I D 2 のこれらの範囲を共に下層媒体 I D 1 に対してマッピングするが、同一のソース媒体の別個の範囲を異なる下層媒体に対してマッピングすることもできる。例えば媒体 I D 3 5 の別個の範囲を別個の下層媒体に対してマッピングする。例えば媒体 I D 3 5 の範囲 0 ~ 2 9 9 を、オフセット 4 0 0 で下層媒体 I D 1 8 に対してマッピングする。これは、媒体 I D 3 5 のブロック 0 ~ 2 9 9 を媒体 I D 1 8 のブロック 4 0 0 ~ 6 9 9 に対してマッピングすることを示す。更に媒体 I D 3 5 の範囲 3 0 0 ~ 4 9 9 を、~ 3 0 0 のオフセットで下層媒体 I D 3 3 に対してマッピングし、媒体 I D 3 5 の範囲 5 0 0 ~ 8 9 9 を、~ 4 0 0 のオフセットで下層媒体 I D 5 に対してマッピングする。これらのエントリは、媒体 I D 3 5 のブロック 3 0 0 ~ 4 9 9 が媒体 I D 3 3 のブロック 0 ~ 1 9 9 に対してマッピングされ、媒体 I D 3 5 のブロック 5 0 0 ~ 8 9 9 が媒体 I D 5 のブロック 1 0 0 ~ 4 9 9 に対してマッピングされることを示す。なお他の実施形態では、媒体は 4 つ以上の範囲に分解してよい。

40

【 0 0 2 9 】

テーブル 3 0 0 の「状態」列は、ブロックのための探索をより効率的に実行できるようにする情報を記録する。状態「Q」は、媒体が静止状態 (q u i e s c e n t) であることを示し、「R」は媒体が登録状態 (r e g i s t e r e d) であることを示し、「U」

50

は媒体が非マスク状態 (unmasked) であることを示す。静止状態では、テーブル 300 内で特定された 1 つ又は 2 つの媒体だけに対して探索を実行する。登録状態では、探索を再帰的に実行する。非マスク状態は、探索を基底媒体において実行するべきか、又は探索を下層媒体においてのみ実行するべきかを決定する。いずれのエントリに関してテーブル 300 には示されていないが、別の状態「X」を用いて、ソース媒体が非マッピング状態であることを特定してよい。この非マッピング状態は、ソース媒体が到達可能なデータを全く含まず、破棄できることを示す。この非マッピング状態はソース媒体の範囲に適用してよい。媒体全体が非マッピング状態である場合、媒体 ID をシーケンス無効化テーブルに入力して最終的に破棄してよい。

【0030】

一実施形態では、媒体が生成されると、この媒体が下層媒体を有する場合はこの媒体は登録状態となり、又はこの媒体が既存の状態を有さない新規のボリュームである場合はこの媒体は静止状態となる。媒体に書き込みが行われると、この媒体の一部は非マスク状態となり得、媒体自体及び下層媒体内にマッピングが存在する。これは、単一の範囲を複数の範囲のエントリに分割することによって実行でき、上記複数の範囲のエントリのうちのいくつかは、初期のマスク状態に保持され、その他は非マスク状態として記録される。

【0031】

更にテーブル 300 内の各エントリは、基底属性を含んでよく、これは媒体の基底を示し、これはこの場合ソース媒体自体を指す。各エントリはまた、オフセットフィールドを含んでよく、これは下層媒体に対してソース媒体にマッピングする際にブロックのアドレスに適用するべきオフセットを特定する。これにより媒体は、下層媒体の開始ブロックから下層媒体の頂部に構成されるだけでなく、下層媒体内の他の位置に対してマッピングできる。テーブル 300 に示すように、媒体 8 は 500 のオフセットを有し、これは媒体 8 のブロック 0 を下層媒体 (媒体 1) のブロック 500 に対してマッピングすることを示す。従って媒体 8 による媒体 1 の探索は、500 のオフセットを、その要求の初期ブロック番号に付加することになる。「オフセット」の行により、媒体を複数の媒体で構成できる。例えば一実施形態では、媒体は「ゴールドマスター (gold master) 」オペレーティングシステム画像及び VM (仮想マシン) 毎のスクラッチスペースからなる。その他の柔軟なマッピングも可能であり、考察の対象となる。

【0032】

各エントリは下層媒体属性も含み、これはソース媒体の下層媒体を示す。(媒体 1 と同様に) 下層媒体がソース媒体を指す場合、これはソース媒体が下層媒体を有さず、全ての探索はソース媒体のみにおいて実行されることを示す。各エントリは安定属性も含み、「Y」 (yes : はい) は媒体が安定であること (又は読み出し専用であること) を示し、「N」 (no : いいえ) は媒体が読み書き可能であることを示す。安定である媒体では、媒体内の所定のブロックに対応するデータは変化しないが、このデータを生成するマッピングは変化し得る。例えば媒体 2 は安定であるが、媒体 2 のブロック 50 は媒体 2 又は媒体 1 に記録してよく、媒体 2 及び媒体 1 はこの順に論理的に探索されるが、この探索は必要に応じて並列に実行してよい。一実施形態では、媒体がこれ自体以外のいずれの媒体によって下層媒体として使用される場合、この媒体は安定となる。

【0033】

ここで図 4 を参照すると、コピーオフロード動作の一実施形態を示す。図 4 に示すコピーオフロード動作は、ボリューム 8 のブロック 400 ~ 599 はボリューム 12 のブロック 850 ~ 1049 にコピーされることになることを特定する。ボリューム 8 及びボリューム 12 の論理表現を、連続データ構造として図 4 に示しているが、ボリューム 8 及びボリューム 12 に対応するデータは、ホストストレージシステムのストレージデバイス全体に亘って、多数の別個の非連続的な位置に配置され得ることを理解されたい。

【0034】

ストレージコントローラは、このコピーオフロード動作を実行するための要求の受信に回答して、ボリューム 8 及びボリューム 12 がどの媒体を指すかを決定できる。一実施形

10

20

30

40

50

態では、ストレージコントローラはボリューム - 媒体間マッピングテーブルに問い合わせを行うことによって、ボリューム 8 及びボリューム 12 のアンカー媒体を決定してよい。この議論の目的に関して、(コピーオフロード動作実行前のボリューム - 媒体間マッピングテーブルの 2 つのエントリを表すボリューム - 媒体間マッピングテーブル 405A に示すように) ボリューム 8 のアンカー媒体が媒体 123 である、及びボリューム 12 のアンカー媒体が媒体 234 であると考えられる。また媒体マッピングテーブル 410A は、コピーオフロード動作実行前の媒体マッピングテーブルの 2 つのエントリを表す。

【0035】

要求されたコピーオフロード動作を実行するために、ストレージコントローラは媒体「x」を生成してよく、媒体 x の下層媒体として媒体 123 を記録してよい。この時点において、媒体 123 は安定となる。値「x」は現在使用中でないいずれの媒体 ID 値であってよい。また、(コピーオフロード動作実行後のボリューム - 媒体間マッピングテーブルの 2 つのエントリを表す) ボリューム - 媒体間マッピングテーブル 405B に示すように、ボリューム 8 は媒体 x を指すように更新してよい。更に媒体「y」を生成してよく、媒体 y の下層媒体として媒体 234 を記録してよい。またボリューム 12 は媒体 y を指すように更新してよい。コピーされる媒体 y の範囲に関してエントリを生成してよく、このエントリは、その下層 (underlying: U) 媒体として媒体 123 を有してよい。他の範囲の媒体 y のための他の全てのエントリは、下層媒体として媒体 234 を有してよい。これらのエントリを媒体マッピングテーブル 410B に示し、これはコピーオフロード動作実行後の媒体マッピングテーブルの一部を表す。なお媒体マッピングテーブル 410A - B のエントリは、図 4 に示しているものに対する追加の情報を含んでよい。

【0036】

上述のステップに従ってコピーオフロード動作を実行することにより、ストレージコントローラは、実際に問題のデータブロックをコピーすることなく、要求されたコピーオフロード動作を実行できる。その代わりにコピーオフロード動作は、ボリューム - 媒体マッピングテーブル及び媒体マッピングテーブルに対して変更を行うだけで遂行される。その結果として、コピーオフロード動作を実装するにあたって即時データ書き込みは実行されず、コピーオフロード動作を最小のリソース利用によって迅速に実行できる。

【0037】

コピーオフロード動作の上述の説明はまた、他のタイプのコピーオフロード動作を実行する異なる実施形態にも応用できる。例えば別の実施形態では、第 1 のボリュームの第 1 の範囲から複数の別個のボリュームに対してコピーオフロード動作を要求してよい。このコピーオフロード動作に関して、媒体「x」及びボリューム 8 に対応する上述の複数のステップを、ソースボリュームに対して 1 回実行してよい。媒体「y」及びボリューム 12 に対応する複数のステップを、このコピーオフロード動作の標的である各対象ボリュームに関して反復してよく、各対象ボリュームに関して新規の媒体を生成する。

【0038】

いくつかの実施形態では、コピーオフロード動作を実行するための上述の複数のステップは、コピーオフロード動作を実行するための要求の受信後すぐに実装するのではなく、バッファリングしてよい。複数のコピーオフロード動作をバッファリングして、これを後にバッチ・モードで実行してよい。更にいくつかの実施形態では、ソースボリュームに対応するコピーオフロード動作の複数のステップを、ソースボリュームを標的とする要求をストレージコントローラが受信するまで遅延させてよい。この時点において、ソースボリュームに対応するコピーオフロード動作の一部を実行してよい。同様に対象ボリュームに対応するコピーオフロード動作の複数のステップを、対象ボリュームを標的とする要求をストレージコントローラが受信するまで遅延させてよい。

【0039】

ここで図 5 を参照すると、コピーオフロード動作の別の実施形態のブロック図が示されている。このコピーオフロード動作では、複数のデータブロックのセットが、あるボリューム内の第 1 の位置から、同じボリューム内の第 2 の位置へとコピーされている。ボリューム

ーム 35 の論理表現を図 5 に示しており、データは位置 200 ~ 499 から位置 1800 ~ 2099 へとコピーされる。この議論の目的に関して、媒体 355 はテーブル 505 A に示すように、ボリューム 35 のアンカー媒体であると考えられる。媒体 355 に関するエントリを媒体マッピングテーブル 510 A に示しており、コピーオフロード動作実行前の媒体マッピングテーブルの 1 つのエントリの代表例である。

【0040】

一実施形態では、このコピーオフロード動作を実行するための要求の受信に回答して、ストレージコントローラは新規の媒体「z」を生成してよい。テーブル 505 B に示すように、媒体 z はボリューム 35 のアンカー媒体として記録してよい。また、媒体 z に関する 3 つの別個のエントリを、媒体マッピングテーブル（その一部分をテーブル 510 B に示す）に追加してよい。媒体 z に関する第 1 のエントリは、データブロックの範囲 0 ~ 1799 のためのものであり、これに対する下層（U）媒体は媒体 355 として記録される。第 1 のエントリのためのオフセットは 0 に設定される。同様に、データブロックの範囲 2100 ~ N のための、媒体 z に関する第 3 のエントリは、第 1 のエントリと同一の属性を有する。第 1 のエントリ及び第 3 のエントリはそれぞれ 0 のオフセットを有し、これは下層媒体（媒体 355）に対するマッピング時に使用される。媒体 z に関する第 2 のエントリは、コピーされたデータが標的とする範囲（1800 ~ 2099）に対応する。第 2 のエントリはまた、その下層媒体として記録された媒体 355 を有する。しかしながらこの第 2 のエントリは ~ 1600 のオフセットを有し、これにより、コピーオフロード動作において特定されたデータに対応する媒体 355 内の正しい位置に対してマッピングできる。

10

20

【0041】

上述の技術を用いることによって、ストレージコントローラは、問題のデータブロックを物理的にコピーすることなく、要求されたコピーオフロード動作を達成できる。寧ろコピーオフロード動作は、データブロックにアクセスせずにボリューム及び媒体マッピングテーブルを操作するだけで実行される。コピーのために要求されるデータが記憶されている実際の物理ストレージ位置は、このコピーオフロード動作中にアクセスされることはない。

【0042】

ここで図 6 に移ると、コピーオフロード動作を実行するための方法 600 の一実施形態が示されている。上述のシステム 100 に統合されている構成部品（例えばストレージコントローラ 110）は一般に、方法 600 に従って動作し得る。更にこの実施形態のステップを順番に示す。しかしながら別の実施形態では、いくつかのステップは図示したものと異なる順序で実行してよく、いくつかのステップは同時に実行してよく、いくつかのステップは他のステップと組み合わせてよく、またいくつかのステップが存在しなくてもよい。

30

【0043】

第 1 のボリュームの第 1 の範囲から第 2 のボリュームの第 2 の範囲へのコピーオフロード動作を実行するための要求を、ストレージコントローラが受信してよい（ブロック 605）。この議論の目的に関して、第 1 のボリュームは第 1 の媒体に関連する（即ち第 1 の媒体を指す）と考えられる。また第 2 のボリュームは第 2 の媒体を指すと考えられる。一実施形態では、これらの関連はボリューム - 媒体マッピングテーブルに問い合わせを行うことによって決定してよい。

40

【0044】

この要求の受信に回答して、第 3 の媒体を生成してよく、媒体マッピングテーブルにおいてこの第 3 の媒体に関する新規のエントリを生成してよい（ブロック 610）。様々な実施形態では、新規の媒体の生成プロセスは、この媒体に関する新規の ID を生成して、媒体マッピングテーブルにおいて上記媒体に関する新規のエントリを生成することを伴う。第 3 の媒体の下層媒体として第 1 の媒体を指定する指示を記憶してよい（ブロック 615）。一実施形態では、ブロック 615 は、媒体マッピングテーブルの新規のエントリに

50

第3の媒体の下層媒体として第1の媒体を記録することによって実装してよい。

【0045】

続いてボリューム - 媒体間マッピングテーブルを更新して、第1のボリュームを第3の媒体に関連させてよい(ブロック620)。換言すると、第3の媒体を第1のボリュームのアンカー媒体として特定してよい。また、第1の媒体を読み出し専用(即ち安定)として特定するインジケータを記憶してよい(ブロック625)。一実施形態では、このインジケータを媒体マッピングテーブル内の対応するエントリに記憶してよい。

【0046】

更に第4の媒体を生成してよく、媒体マッピングテーブル内でこの第4の媒体に関する新規のエントリを生成してよい(ブロック630)。第2の媒体を、第4の媒体の下層媒体として指定してよい(ブロック635)。また、第1の媒体の第1の範囲が第4の媒体の第2の範囲の下層となるよう指定するインジケータを記憶してよい(ブロック640)。第4の媒体の第2の範囲に関して、別個の範囲のエントリを媒体マッピングテーブル内に生成してよく、この別個の範囲のエントリは、第4の媒体の第2の範囲を第1の媒体の第1の範囲に対してマッピングできる。一実施形態では、オフセット値をこの範囲のエントリに含むことにより、第4の媒体の第2の範囲を第1の媒体の第1の範囲に対して整列させる方法を特定してよい。第2の媒体は、第4の媒体の他の全ての範囲に関して下層媒体として保持され得る。

10

【0047】

更に第4の媒体を、第2のボリュームのアンカー媒体として特定してよい(ブロック645)。また更に、第2の媒体を読み出し専用と特定するインジケータを記憶してよい(ブロック650)。ブロック650の後、方法600を終了してよい。上述の方法を用いて、いずれの数のコピーオフロード動作を並行して実行してよい。なお上述のステップのうちいくつかを、コピーオフセット動作が作用するボリュームの領域を標的とする読み出し又は書き込み動作が受信されるまで遅延させてよい。

20

【0048】

ここで図7を参照すると、コピーオフロード動作のための別の方法700の一実施形態が示されている。上述のシステム100に統合されている構成部品(例えばストレージコントローラ110)は一般に、方法700に従って動作し得る。更にこの実施形態のステップを順番に示す。しかしながら別の実施形態では、いくつかのステップは図示したものと異なる順序で実行してよく、いくつかのステップは同時に実行してよく、いくつかのステップは他のステップと組み合わせるとよく、またいくつかのステップが存在しなくてもよい。

30

【0049】

コピーオフロード動作を実行するための要求をストレージコントローラが受信してよい(ブロック705)。この議論の目的に関して、コピーオフロード動作の目標は、第1のボリュームの第1の範囲から第2のボリュームの第2の範囲にコピーされることになるデータであると考えられる。しかしながら他の動作は、異なるソース及び/又は対象を有してよい。要求の受信に回答して、受信したコピーオフロード要求に対応する動作は、即座に実行するのではなくバッファリングしてよい(ブロック710)。様々な実施形態では、複数の動作を別個にバッファリングしてよい。例えば第1のボリュームの第1の範囲を後続の要求が標的とした時に、第1のボリュームに対応する動作を生成し、バッファリングし、実行準備状態としてよい。同様に、第2のボリュームの第2の範囲が標的とされた時に、第2のボリュームに対応する動作を生成し、バッファリングし、実行準備状態とする。

40

【0050】

ブロック710の後、ストレージコントローラが、過去に受信したコピーオフロード動作が作用する領域を標的とする読み出し又は書き込み要求を受信したかどうかを決定してよい(条件ブロック715)。作用される位置を標的とする要求が受信されていない場合(条件ブロック715の「いいえ(No)」分岐)、ストレージコントローラは、1

50

つ又は複数のバッファリングされたコピーオフロード動作の実行を防止できる（ブロック720）。ブロック720の後、別のコピーオフロード要求を受信したかどうかを決定してよい（ブロック725）。コピーオフロード要求を受信されていた場合（条件ブロック725の「はい（Yes）」分岐）、方法700はブロック710に戻り、受信したコピーオフロード要求を実行できる。なお、ストレージコントローラは、コピーオフロード要求を（ブロック725が表す時点だけではなく）いずれの時点において受信してよく、この場合方法700はこれに従ってブロック710にジャンプして、受信したコピーオフロード要求をバッファリングしてよい。コピーオフロード要求を受信されていなかった場合（条件725の「いいえ」分岐）、方法700はブロック715に戻り、バッファリングされたコピーオフロード要求に対応する領域を標的とする要求を受信したかどうかを決定してよい。

10

【0051】

作用される位置を標的とする要求を受信されていた場合（条件ブロック715の「はい」分岐）、ストレージコントローラによって対応するコピーオフロード動作を実行してよい（ブロック730）。バッファリングされたコピーオフロード動作は、ソースボリューム及び対象ボリュームを標的としてよく、受信された要求がこれらのボリュームのうちの1つのみを標的とする場合、上記作用されるボリュームを標的とするコピーオフロード動作のみを実行してよい。コピーオフロード動作のその他の部分はバッファリングされたままであってよく、（対応するボリューム内における上記部分の位置が標的とされた場合、又は処理リソースがアイドル状態であり、使用できる状態である場合に）後の時点で実行してよい。ブロック730の後、方法700はブロック725に戻り、別のコピーオフロード要求を受信したかどうかを決定してよい。

20

【0052】

作用される領域（又はボリューム）を標的とする後続の要求を受信されるまで、コピーオフロード動作の実行を待機することによって、ストレージシステムの処理リソースを、他のタスクを実行できるよう自由にすることができる。コピーオフロード動作のバッファリングはまた、追加の媒体の生成を、それが実際に必要となるまで防止することによって補助できる。ストレージシステムに不要な負荷を与えることなく、複数のコピーオフロード動作を受信してバッファリングできる。また、ストレージシステムのリソースが利用可能な期間中、ストレージコントローラはアイドル状態の処理性能を用いて、多数のバッファリングされたコピーオフロード動作を実行してよい。このようにして、ストレージシステムが実行している他のタスクに干渉することなく、コピーオフロード動作を実行できる。従っていくつかの実施形態では、ブロック710が指示するように、全ての受信したコピーオフロード動作をバッファリングするのではなく、ストレージコントローラは、受信したコピーオフロード動作がケースバイケースでバッファリングされているかどうかを決定してよい。この決定は、少なくともストレージシステムの現在の実行条件（例えば処理負荷、ストレージ利用、保留中の要求の個数）に基づいてよい。他の実施形態では、受信したコピーオフロード動作を自動的にバッファリングしてよく、またバッファリングされるコピーオフロード動作の個数が閾値を超えると、ストレージコントローラは、バッチ・モードで複数のコピーオフロード動作を実行してよい。これらの実施形態では、作用される領域を標的とするデータ要求を受信すると、対応するバッファリングされたコピーオフロード動作を、その他のコピーオフロード動作がバッファリングされたままの状態で行ってよい。

30

40

【0053】

なお上述の実施形態はソフトウェアを備えてよい。このような実施形態では、本方法及び/又は機構を実装するプログラム命令を、コンピュータ可読媒体上で搬送又は記憶してよい。プログラム命令を記憶するよう構成された多数のタイプの媒体が利用可能であり、これにはハードディスク、フロッピーディスク、CD-ROM、DVD、フラッシュメモリ、プログラマブルROM（PROM）、ランダムアクセスメモリ（RAM）、他の様々な形態の揮発性又は不揮発性ストレージが含まれる。

50

【0054】

様々な実施形態において、本明細書に記載の方法及び機構の1つ又は複数の部分は、クラウドコンピューティング環境の一部を形成してよい。このような実施形態では、1つ又は複数の様々なモデルに従って、必要に応じてインターネットを介してリソースを提供してよい。このようなモデルは、インフラストラクチャ・アズ・ア・サービス (Infrastructure as a Service: IaaS)、プラットフォーム・アズ・ア・サービス (Platform as a Service: PaaS)、ソフトウェア・アズ・ア・サービス (Software as a Service: SaaS) を含んでよい。IaaSでは、コンピュータインフラストラクチャを必要に応じて提供する。このような場合、一般にサービスプロバイダがコンピュータ設備を所有し、動作させる。PaaSモデルでは、サービスプロバイダが、開発者がソフトウェアによる解決法を開発するために使用するソフトウェアツール及びその下層設備を必要に応じて提供し、そのホストとなってよい。SaaSは典型的には、サービスに応じてソフトウェアを認可するサービスプロバイダを含む。このサービスプロバイダはソフトウェアのホストとなってよく、又はソフトウェアを所定の期間だけカスタマに対して配備してよい。上述のモデルの多数の組み合わせが可能であり、考察の対象となる。

10

【0055】

以上の実施形態についてかなり詳細に説明したが、上述の開示を完全に理解すれば、当業者には多数の変形例及び修正例が明らかとなるであろう。以下の請求項は、これらの変形例及び修正例の全てを包含するものと解釈されることを意図したものである。

20

【図1】

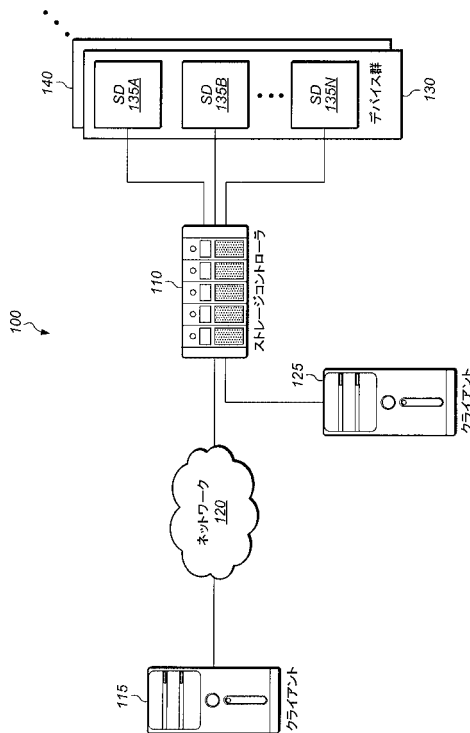


FIG. 1

【図2】

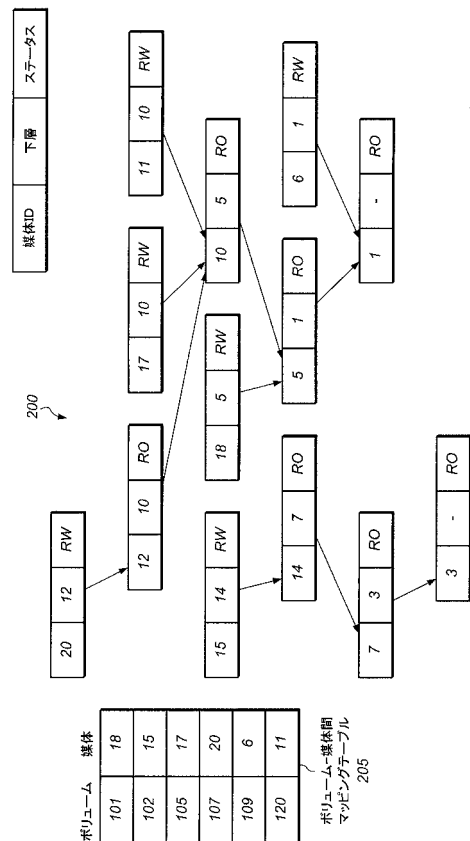


FIG. 2

【図3】

媒体ID	範囲	状態	基底	オフセット	下階	安定
1	0-999	Q	1	0	1	Y
2	0-99	QU	2	0	1	Y
2	100-999	Q	2	0	1	Y
5	0-999	RU	5	0	2	N
8	0-499	R	8	500	1	N
10	0-999	QU	10	0	1	Y
14	0-999	RU	14	0	10	Y
18	0-999	RU	18	0	14	N
25	0-999	RU	25	0	14	Y
33	0-999	RU	33	0	25	N
35	0-299	RU	35	400	18	N
35	300-499	RU	35	-300	33	Y
35	500-899	RU	35	-400	5	N

Q - 静止状態; R - 登録状態; U - 非マスク状態

FIG. 3

【図4】

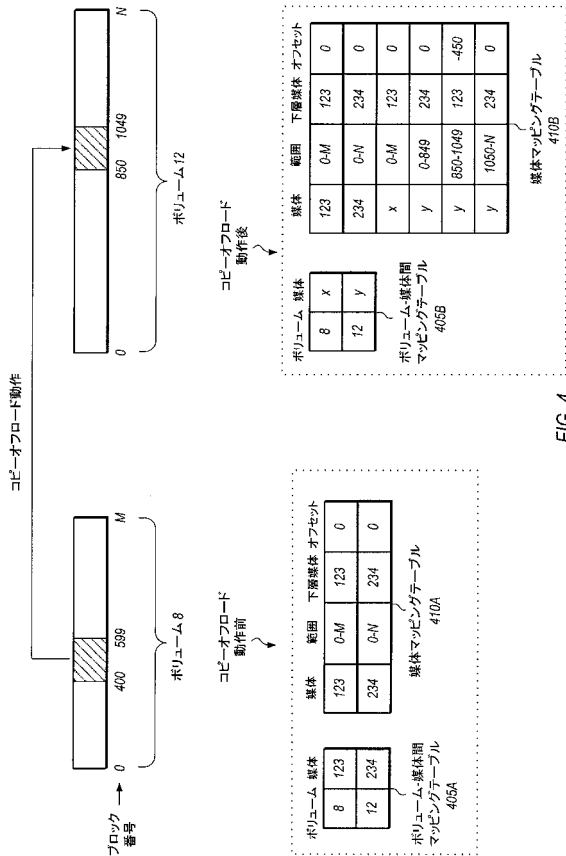


FIG. 4

【図5】

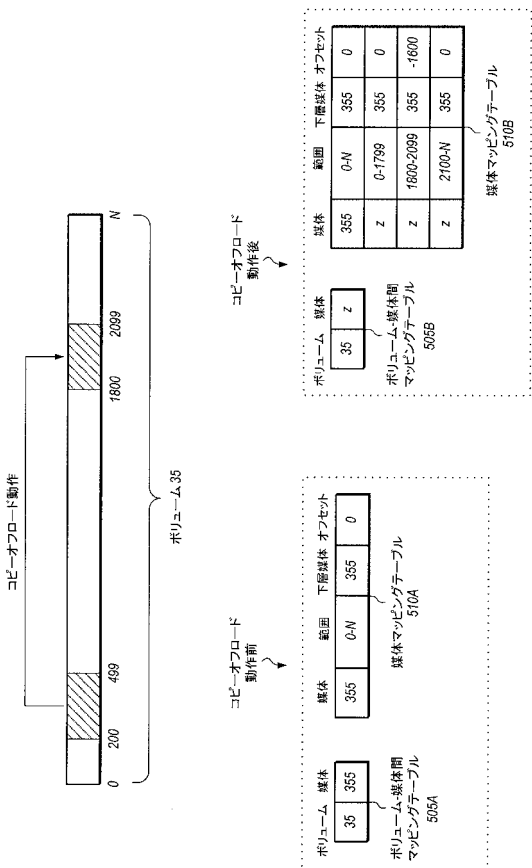


FIG. 5

【図6】

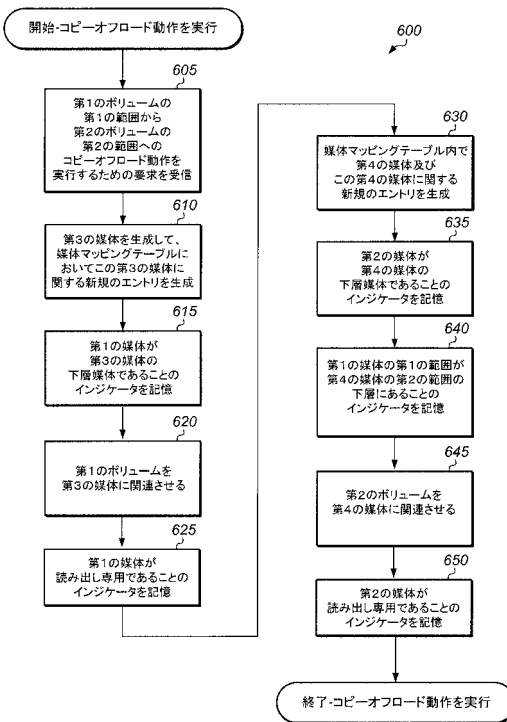


FIG. 6

【 図 7 】

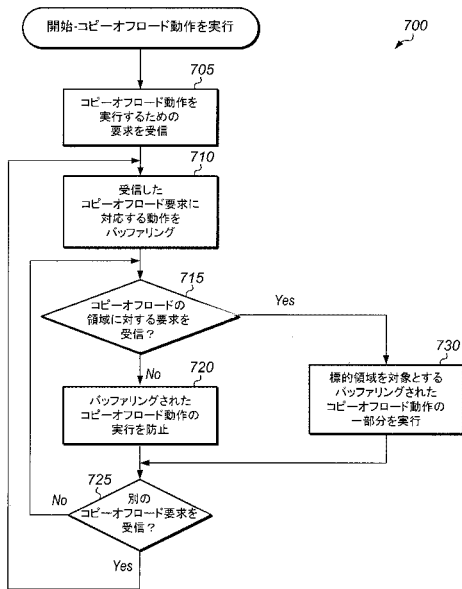


FIG. 7

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No PCT/US2014/010690

A. CLASSIFICATION OF SUBJECT MATTER INV. G06F3/06 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) G06F		
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		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	US 8 290 911 B1 (JANAKIRAMAN VISWESVARAN [US] ET AL) 16 October 2012 (2012-10-16) figures 3-6 column 7, line 17 - column 10, line 33 -----	1-22
Y	US 2006/174074 A1 (BANIKAZEMI MOHAMMAD [US] ET AL) 3 August 2006 (2006-08-03) figures 2,4 paragraph [0024] paragraph [0026] -----	1-22
A	US 2012/330903 A1 (PERIYAGARAM SUBRAMANIAM [US] ET AL) 27 December 2012 (2012-12-27) paragraph [0075] - paragraph [0083] figures 8-10 -----	1-22
<input 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
11 March 2014		24/03/2014
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 Andlauer, J

1

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2014/010690

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 8290911	B1	16-10-2012	NONE

US 2006174074	A1	03-08-2006	CA 2593289 A1 10-08-2006
			CN 101120305 A 06-02-2008
			EP 1853992 A2 14-11-2007
			JP 2008529187 A 31-07-2008
			US 2006174074 A1 03-08-2006
			WO 2006083327 A2 10-08-2006

US 2012330903	A1	27-12-2012	US 2012330903 A1 27-12-2012
			WO 2012177318 A1 27-12-2012

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, 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, LT, 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

(特許庁注：以下のものは登録商標)

1 . E T H E R N E T

- (72)発明者 ミラー, イーサン
アメリカ合衆国・95060・カリフォルニア州・サンタクルーズ・カルカー ドライブ・203
- (72)発明者 ヘイズ, ジョン
アメリカ合衆国・94041・カリフォルニア州・マウンテン ビュー・ハイスクール ウェイ・800・ナンバー・330
- (72)発明者 サンドヴィグ, ケイリー
アメリカ合衆国・94303・カリフォルニア州・パロ アルト・ドノホー ストリート・284
- (72)発明者 ゴールデン, クリストファー
アメリカ合衆国・94041・カリフォルニア州・マウンテン ビュー・ハイスクール ウェイ・800・ナンバー・229
- (72)発明者 カオ, ジャンティン
アメリカ合衆国・94043・カリフォルニア州・マウンテン ビュー・サン カルロス アヴェニュー・791
- (72)発明者 イノゼムツェフ, グリゴリ
アメリカ合衆国・94085・カリフォルニア州・サニーベイル・エスカロン アヴェニュー・1000・ナンバー3020

Fターム(参考) 5B005 MM51