

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6737981号
(P6737981)

(45) 発行日 令和2年8月12日 (2020.8.12)

(24) 登録日 令和2年7月21日 (2020.7.21)

(51) Int. Cl.

F I

G 0 6 F 13/28 (2006.01)

G 0 6 F 13/28 3 1 0 A

G 0 6 F 3/06 (2006.01)

G 0 6 F 13/28 3 1 0 E

G 0 6 F 3/06 3 0 1 Z

請求項の数 9 (全 12 頁)

(21) 出願番号 特願2017-533309 (P2017-533309)
 (86) (22) 出願日 平成27年12月22日 (2015.12.22)
 (65) 公表番号 特表2018-504689 (P2018-504689A)
 (43) 公表日 平成30年2月15日 (2018.2.15)
 (86) 国際出願番号 PCT/EP2015/080915
 (87) 国際公開番号 W02016/110410
 (87) 国際公開日 平成28年7月14日 (2016.7.14)
 審査請求日 平成30年8月17日 (2018.8.17)
 (31) 優先権主張番号 14/589,758
 (32) 優先日 平成27年1月5日 (2015.1.5)
 (33) 優先権主張国・地域又は機関
 米国 (US)

(73) 特許権者 390009531
 インターナショナル・ビジネス・マシーンズ・コーポレーション
 INTERNATIONAL BUSINESS MACHINES CORPORATION
 アメリカ合衆国10504 ニューヨーク州アーモンク ニュー オーチャードロード
 New Orchard Road, Armonk, New York 10504, United States of America
 (74) 代理人 100108501
 弁理士 上野 剛史

最終頁に続く

(54) 【発明の名称】 リモート・ダイレクト・メモリ・アクセス操作を実行するデータ転送方法、システム、およびプログラム

(57) 【特許請求の範囲】

【請求項 1】

メモリ・デバイスによる、ネットワーク・ファイル・システム (NFS) のリモート・ダイレクト・メモリ・アクセス (RDMA) 操作における効率的なデータ転送のための方法であって、

前記 RDMA 操作を行うことに応じて、データ自体に触れることなく、前記データのファイル・プロトコル・ヘッダをブロック・プロトコル・ヘッダに置き換えること

を含み、前記ブロック・プロトコル・ヘッダは、スモール・コンピューティング・システム・インターフェース (SCSI) RDMA プロトコル (SRP) を用いて、中央処理ユニット (CPU) による外部からの操作なしに、トランスポート層を介して直接データ経路上でのソースからターゲットへの転送を可能にし、それによって、前記 RDMA のデータ転送は低レイテンシのネットワークングを可能にする、前記方法。

【請求項 2】

前記ファイル・プロトコル・ヘッダを置き換えることと併せて、前記データのチャンクへの少なくとも 1 つのポインタを前記ブロック・プロトコル・ヘッダに関連付けて設定することをさらに含む、請求項 1 に記載の方法。

【請求項 3】

前記ファイル・プロトコル・ヘッダが表すファイルのディスク割り当てスキームを解決することをさらに含む、請求項 1 又は 2 に記載の方法。

【請求項 4】

前記ブロック・プロトコル・ヘッダを生成することをさらに含む、請求項 3 に記載の方法。

【請求項 5】

前記ファイル・プロトコル・ヘッダを構文解析することをさらに含む、請求項 1 ~ 4 のいずれか 1 項に記載の方法。

【請求項 6】

入ってくるファイル・プロトコル要求を受け取ることをさらに含む、請求項 1 ~ 5 のいずれか 1 項に記載の方法。

【請求項 7】

前記置き換えにおいて、前記ファイル・プロトコル・ヘッダ中のファイル ID、オフセット及びサイズ・ヘッダ情報がストレージのブロック番号を含むヘッダに置き換えられる、請求項 1 ~ 6 のいずれか 1 項に記載の方法。

10

【請求項 8】

ネットワーク・ファイル・システム (NFS) のリモート・ダイレクト・メモリ・アクセス (RDMA) 操作における効率的なデータ転送のためのシステムであって、

メモリ・デバイスを含み、前記メモリ・デバイスは、請求項 1 ~ 7 のいずれか 1 項に記載の方法を実行することを可能にする、前記システム。

【請求項 9】

メモリ・デバイスによる、ネットワーク・ファイル・システム (NFS) のリモート・ダイレクト・メモリ・アクセス (RDMA) 操作をコンピュータ・システムにおいて効率的にデータ転送するためのコンピュータ・プログラムであって、前記コンピュータ・システムに、請求項 1 ~ 7 のいずれか 1 項に記載の方法の各ステップを実行させる、前記コンピュータ・プログラム。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に、コンピューティング・システムに関し、より特定的には、リモート・ダイレクト・メモリ・アクセス (Remote Direct Memory Access、RDMA) 操作を実施するクラスタ化されたコンピューティング・システムにおける効率的なデータ転送のための種々の実施形態に関する。

30

【背景技術】

【0002】

最新技術を用いる今日では、大量のデータをディスク及びフラッシュ・ドライブ上に格納可能であり、これらのドライブは、コンピューティング・ストレージ・ネットワークのような大規模ストレージ環境内に単独のエンティティとして又はより広い構成の一部として存在し得る。今日の情報ベースの社会における膨大な量のデータが増大し続けるにつれて、コンピューティング及びコンピューティング・ストレージ・ネットワークのサイズ及び複雑さも増大し続ける。今日、幾つかの個々のコンピュータのグループ又はクラスタが、データ・ストレージ及び転送を容易にするのは珍しいことではない。

【発明の概要】

40

【発明が解決しようとする課題】

【0003】

リモート・ダイレクト・メモリ・アクセス (RDMA) 操作における効率的なデータ転送のための方法、システム及びコンピュータ・プログラムを提供する

【課題を解決するための手段】

【0004】

現在のファイル・ストレージ・アレイは、入力/出力 (I/O) 要求を、分離層において間接的な方法で処理する。その際、I/O 要求のレイテンシがもたらされ、性能が低下することがある。新しいフラッシュ技術の実装には、ファイル・ストレージが超低レイテンシで何百万もの I/O 作業を収容できる、対応する必要性が付随する。

50

【 0 0 0 5 】

リモート・ダイレクト・メモリ・アクセス (R D M A) は、いずれのコンピュータのオペレーティング・システム (O S) も関与しない、1つのコンピュータのメモリから別のコンピュータのメモリへの直接メモリ・アクセス操作である。R D M A データ転送は、高スループット、低レイテンシのネットワークングを可能にし、そのことは、例えば、ストレージ・ソリューションにおいて特に有利である。

【 0 0 0 6 】

ファイル・ストレージ・アレイは、一般に、割り込み駆動型モデルを使用する。このモデルには、ボトルネック及びレイテンシの増大といった潜在的な問題が付随する。ストレージ・アレイ・アーキテクチャも、膨大な量のコンテキスト・スイッチを使用する。これらのスイッチング操作は、単一のファイル操作要求の実行経路において、異なるユーザとカーネル・スレッドとの間で行われる。これはレイテンシをさらに大きくする傾向があり、性能をさらに低下させる。

【 0 0 0 7 】

従って、上記に鑑みて、コンピューティング・ストレージ環境におけるレイテンシを危うくすることなく、大量のデータ及び/又は多くの I / O 操作を行うことができる機構に対する必要性が存在する。

【 0 0 0 8 】

この必要性に対処するために、メモリ・デバイスによる、リモート・ダイレクト・メモリ・アクセス (R D M A) 操作における効率的なデータ転送のための種々の実施形態が提供される。一実施形態において、単なる例として、メモリ・デバイスによる、R D M A 操作における効率的なデータ転送のための方法が提供される。データのファイル・プロトコル・ヘッダがブロック・プロトコル・ヘッダに置き換えられる。ブロック・プロトコル・ヘッダは、中央処理ユニット (C P U) による外部からの操作なしに、ソースからターゲットへの、トランスポート層を介して直接データ経路上での転送を可能にする。

【 0 0 0 9 】

他のシステム及びコンピュータ・プログラムの実施形態が提供され、関連した利点を提供する。

【 0 0 1 0 】

本発明の利点を容易に理解するために、添付図面に示される特定の実施形態を参照することによって、上記に簡単に説明される本発明のさらに具体的な説明が与えられる。これらの図面は本発明の典型的な実施形態を示すものに過ぎず、従って、その範囲を限定していませんと考えるべきではないことを理解した上で、添付図面を用いて付加的な具体性及び詳細さをもって本発明を記載し、説明する。

【図面の簡単な説明】

【 0 0 1 1 】

【図 1】本発明の態様を実現することができる、リモート・ダイレクト・メモリ・アクセス (R D M A) 操作において効率的なデータ転送を行うための例示的なハードウェア構造を示すブロック図である。

【図 2】本発明の態様を実現することができる、ファイルからブロックへのヘッダ変換 (file to block header transformation) のための例示的なファイルのブロック図である。

【図 3】グローバル・パラレル・ファイル・システム (G l o b a l P a r a l l e l F i l e S y s t e m 、 G P F S) のようなファイル・システムにおける従来のデータ転送のブロック図である。

【図 4】本発明の態様を実現することができる、ファイル・システムにおける直接データ経路データ転送のブロック図である。

【図 5】ここで再び本発明の態様を実現することができる、ヘッダ変換及びダイレクトデータ経路データ転送のための例示的方法のフローチャート図である。

【図 6】ここで再び本発明の態様を実現することができる、ヘッダ変換及び直接データ経

10

20

30

40

50

路データ転送のための付加的な例示的方法の付加的なフローチャート図である。

【発明を実施するための形態】

【0012】

前述のように、現在のファイル・ストレージ・アレイは、入力／出力（I／O）要求を、分離層において間接的な方法で処理する。その際に、I／O要求におけるレイテンシが結果として生じ、性能が低下することがある。新しいフラッシュ技術の実装には、ファイル・ストレージが超低（ultra low）レイテンシで何百万ものI／O操作を収容できる対応する必要性が付随する。

【0013】

リモート・ダイレクト・メモリ・アクセス（RDMA）は、いずれのコンピュータのオペレーティング・システム（OS）も関与することのない、1つのコンピュータのメモリから別のコンピュータのメモリへの直接メモリ・アクセス操作である。RDMAデータ転送は、高スループット、低レイテンシのネットワークングを可能にし、そのことは、例えば、ストレージ・ソリューションにおいて特に有利である。

【0014】

ファイル・ストレージ・アレイは、一般に、割り込み駆動型モデルを使用する。このモデルは、ボトルネック及びレイテンシの増大といった付随する潜在的な問題を有する。ストレージ・アレイ・アーキテクチャもまた、膨大な量のコンテキスト・スイッチを使用する。これらのスイッチング操作は、単一のファイル操作要求の実行経路において異なるユーザとカーネル・スレッドとの間で行われる。これはレイテンシをさらに大きくする傾向があり、性能をさらに低下させる。

【0015】

従って、上記に鑑みて、コンピューティング・ストレージ環境におけるレイテンシを危うくすることなく、大量のデータ及び／又は多くのI／O操作を行うことができる機構に対する必要性が存在する。

【0016】

示される実施形態は、ファイル・ストレージ・アレイに対して直接データ・フローの革新的な機構を導入し、ファイル・ストレージにおける超低レイテンシを達成することにより、この必要性に対処する。これは、例えば、RDMA読み取り及びRDMA書き込み操作におけるファイル・プロトコルからブロック・プロトコルへのヘッダ変換を用いることによって実施される。

【0017】

示される実施形態の機構は、何百万もの秒当たりの入力／出力操作（Input / Output Operations Per Second、IOPS）の容易化及び処理を可能にし、超低レイテンシを実現する。この機構により、新しい技術のいかなる性能特性の損失もなしに、ファイル操作の新しいフラッシュ技術の使用が可能になる。

【0018】

一実施形態において、これらの利点は、トランスポート層がアプリケーション層に関する知識をもって(knowledgeable)おり、逆もまた同様である、新しいモノリシック・モデルの使用により達成される。この層間の共有知識の結果、性能向上のための抽象化(abstraction for better performance)は行われない。

【0019】

図示される実施形態において、中央処理ユニット、すなわちCPUは、データのコピーを処理しない。どちらかと言えば、データは、入力エンドポイント及び出力エンドポイントの両方でRDMAを用いて転送される。

【0020】

ここで図1を参照すると、コンピューティング・システム環境の例示的なアーキテクチャ10が示される。一実施形態において、アーキテクチャ10は、本発明の機構を実施するためのシステムの少なくとも部分として実装することができる。コンピュータ・システム10は、通信ポート18及びメモリ・デバイス16に接続された中央処理ユニット（C

10

20

30

40

50

PU) 12を含む。通信ポート18は、通信ネットワーク20と通信する。通信ネットワーク20及びストレージ・ネットワークは、サーバ(ホスト)24及びストレージ・デバイス14を含むことができるストレージ・システムと通信するように構成することができる。ストレージ・システムは、ハード・ディスク・ドライブ(HDD)デバイス、ソリッド・ステート・デバイス(SSD)等を含むことができ、これらは、redundant array of independent disks(RAID)で構成することができる。通信ポート18、通信ネットワーク20、及び簡潔にするために図示されていない当業者には公知の他のコンポーネントは、ファイバ・チャネル・ケーブル、ファイバ・チャネル・ポート、ホスト・バス・アダプタ(HBA)、コンバージド・ネットワーク・アダプタ(Converged Network Adapter、CNA)、ネットワーク・スイッチ及びスイッチング・コンポーネント、並びに当業者には公知の同様の通信機構のようなハードウェア・コンポーネントを含むことができる。さらに説明されるようなこれらのコンポーネントの1つ又は複数を用いて、図示される実施形態の種々の態様を実現することができる。

【0021】

以下に記載される動作は、システム10又は他の場所に配置され、独立して及び/又は他のCPUデバイス12と共に動作する複数のメモリ・デバイス16を有するストレージ・デバイス14上で実行することができる。本明細書で提示されるようなメモリ・デバイス16は、電氣的消去可能プログラム可能読み出し専用メモリ(EEPROM)のようなメモリ、RDMA操作の実行を任されるデバイス(RDMAカードなど)、又は関連するデバイスのホストを含むことができる。メモリ・デバイス16及びストレージ・デバイス14は、信号搬送媒体を介してCPU12に接続される。さらに、CPU12は、通信ポート18を通じて、取り付けられた複数の付加的なコンピュータ・ホスト・システム24を有する通信ネットワーク20に接続される。さらに、メモリ・デバイス16及びCPU12は、コンピュータ・システム10の各コンポーネント内に埋め込むこと及び含ませることができる。各ストレージ・システムはまた、別個のメモリ・デバイス16及び/又はCPU12と共に、又は別個のメモリ・デバイス16及び/又はCPU12として動作する別個の及び/又は異なるメモリ・デバイス16及び/又はCPU12を含むこともできる。

【0022】

アーキテクチャ10は、コンピュータのクラスタの部分を表すと考えることができ、そこで、CPU12は、別のコンピュータ・システム22内のCPU28及び大容量記憶装置30と通信する。図示のように、メモリ16と26の間のデータ転送を容易にする直接データ経路32と共に、トランスポート層34の一部及びアプリケーション層36も、アーキテクチャ10の部分として示される。

【0023】

当業者であれば、クラスタ化されたコンピューティング環境内の可能なコンピューティング・コンポーネントの全体を必ずしも示していないことを認識するが、ブロック・アーキテクチャは、示される実施形態に関連するような機能を示すように意図される。例えば、アプリケーション層36の部分は、トランスポート層34の部分と通信するように示される。これは、ブロック図の観点から、トランスポート層34がアプリケーション層36に関する知識をもつようにされ、逆もまた同様である、前述のモノリシック・モデルを示すように意図される。メモリ16及び26、並びに直接データ経路32と関連した特定の機能をさらに説明する。

【0024】

次の図2は、ブロック図の形態で、クライアント・ホスト内のファイル・プロトコル・デバイス202から、示されるプロセスを通して、ブロック・プロトコル・デバイス206になる例示的な変換200を示す。より具体的には、当業者であれば理解するように、変換200は、組み込まれたユーザ空間208及びファイル・システム210を有するファイラ204を示す。

【0025】

ファイル・システム210において、ファイル・プロトコル・デバイス202は、RDMAファイル入力/出力(I/O)機能212によって変換プロセスを入力する。プロセスにおける後のステップとして、ファイルには、ファイルからブロックへのヘッダ変換操作214が施され、そこで適切なブロック・プロトコル・ヘッダが生成され、ファイル・プログラム・ヘッダ空間において置換される。プロセスにおける後のステップにおいて、RDMAブロックI/O機能216は、ファイル・プロトコル・デバイスからブロック・プロトコル・デバイス206への変換を完了するように動作する。

【0026】

次の図3は、同じくブロック図の形態で、この後に続き、後で説明する図4との比較のために、グローバル・パラレル・ファイル・システム(Global Parallel File System、GPFS)のようなファイル・システム(しかし、任意の一般的なファイル・システムを含むことができる)を通じた、I/O要求300の形態の例示的な従来のデータ移動を示す。図3及び図4は、GPFSに関連する機能を説明することができるが、当業者であれば、示される実施形態の種々の機構は、任意のファイル・システムの実装に適用可能であり、特定の状況に応じて変化し得ることを理解するであろう。

10

【0027】

要求300は、ユーザ空間302、カーネル空間304、及びファイルシステム・コア306との関連で示される。最初に、要求300を、示されるようなネットワーク・ファイル・システム(Network File System、NFS)のリモート・データ・メモリ・アクセス(RDMA)操作308として受け取り、この要求300は、関連するNFSサーバ310に対して作製され/送られ、次に仮想ファイル・システム(Virtual File System、VFS)312に送られ、次にカーネル拡張314として提供される。

20

【0028】

この時点で、メールボックス/メールボックス・メッセージ処理により、ファイルシステム・コア306による前述した適切なカーネル拡張314の通信が可能になる。メッセージ処理は、ファイルシステム・コア316と拡張314との間で送られ、その時点で、ブロック・デバイス316が構築される。次に、ブロック・デバイス316は、1つ又は複数のSmall Computing Systems Interface(SCSI)層(SCSIアプリケーション層などの)、SCSIトランスポート・プロトコル層(STPL)、又はSCSI相互接続層(SIL)を通じて通信される。次に、ブロック・デバイスは、例えば、SCSI RDMAプロトコル(SRP) RDMA操作320を用いて、別のコンピュータに送られる。

30

【0029】

次の図4は、比較して、同じくブロック図の形態で、本発明の種々の機構が実装されるファイル・システムにおける例示的なデータ移動を示す。図4において、ユーザ空間402及びカーネル空間404も再び示される。要求400を、ユーザ空間402内に受け取り、この要求400は、関連のあるNFSサーバ410により実施される操作と共に、NFS RDMA操作408として開始される。

40

【0030】

VFS312及びカーネル拡張314の関与といった、以前に図3に示される種々の付加的なプロセス・ステップとは対照的に、図4に示されるプロセス・ステップは、最新式である。ここで、直接データ経路により、NFS RDMA操作408が、SCSIイニシエータ412の容易化を伴って実施されるSRP RDMA操作と接続される。次に、SRP RDMAプロトコル414を用いて、データを別のコンピュータに送ることができる。

【0031】

ここで図5を参照すると、方法500として前述したようなヘッダ変換機能を通じて、

50

効率的なデータ転送を実施するための例示的な方法である。方法 500 が開始し（ステップ 502）、ファイル・プロトコル・ヘッダをブロック・プロトコル・ヘッダに置き換え、CPU 操作なしに、直接データ経路転送を可能にする（ステップ 504）。ヘッダ情報だけを置換し、基礎となるデータ自体には触れないことに留意することが重要である。次に、方法 500 は終了する（ステップ 506）。

【0032】

次の図 6 における方法 600 は、より詳細に前述した例示的なヘッダ変換機能を示す。方法 600 が開始し（ステップ 602）、RDMA ハードウェアを用いて、入ってくるファイル・プロトコル要求を受け取る（ステップ 604）。次に、入ってくる要求に従って、ファイル・プロトコル・ヘッダを構文解析する（ステップ 606）。ファイルのディスク割り当てスキーム（すなわち、ファイル・システム論理演算）を解決する（ステップ 608）。この結果、例えば、一対多マッピングがもたらされ得る。

10

【0033】

次のステップ 610 において、以前のファイル・システムにおける情報に基づいて、適切なブロック・プロトコル・ヘッダ情報を生成する。要求関連データ・バッファへの関連したデータ・チャンクへのポインタを設定する（ステップ 612）。次に、必要とされる操作に従って、ブロック関連操作が実施され、RDMA ハードウェアを用いて、ブロック・プロトコル・ヘッダ及び関連したデータ・チャンク・ポインタを生成する（ステップ 614）。

【0034】

20

後の時点で、次に、逆変換がファイルに実施され、結果のファイル・プロトコル・ヘッダ及び随意的なデータ（読み取り操作における）を生成する（ステップ 616）。次に、方法 600 は終了する（ステップ 618）。

【0035】

NFS 書き込み操作との関連で前述した方法 600 のステップの次の例を考える。NFS 書き込み操作ヘッダは、書き込まれるべきデータと一緒に、ファイル ID、オフセット、及びサイズを含む。図 6 を参照すると、ステップ 610 において、ファイル・システム論理は、データ自体に触れることなく、この前述のファイル ID、オフセット及びサイズ・ヘッダ情報を、基礎となるストレージのブロック番号を含む関連するヘッダに置き換える。

30

【0036】

前述のように、示される実施形態の機構は、適用可能なトランスポート層がアプリケーションを認識しており、逆もまた同様であるモノリシック・ソリューションを用いることにより可能にされる。言い換えれば、例えば、トランスポート層及びアプリケーション層の両方とも、互い他方の本質的知識 (intrinsic knowledge) を有し、データ移動操作を実施するための RDMA ベースのプロトコル及びハードウェアを使用する。

【0037】

示される実施形態の態様の 1 つにおいて、RDMA ハードウェア及び関連したプロトコルのセットのみが、配線との間でデータ・バッファに出入りする直接伝送を可能にすることに再度留意されたい。これは、操作又は移動（例えば、カーネル拡張からユーザ空間への移動）なしに、データに触れないことを可能にする。対照的に、示される実施形態の機構は、ヘッダ情報のみを操作する、これは、RDMA データが生来位置合わせされているために可能になり、従って、変換プロセス中に行われるブロック・ストレージの位置合わせ制約を処理する。

40

【0038】

本発明は、システム、方法、及び / 又はコンピュータ・プログラム製品とすることができる。コンピュータ・プログラム製品は、プロセッサに本発明の態様を実施させるためのコンピュータ可読プログラム命令をその上に有するコンピュータ可読ストレージ媒体（単数又は複数）を含むことができる。

【0039】

50

コンピュータ可読ストレージ媒体は、命令実行デバイスにより使用される命令を保持及び格納することができる有形デバイスとすることができる。コンピュータ可読ストレージ媒体は、例えば、これらに限定されるものではないが、電子記憶装置、磁気記憶装置、光学記憶装置、電磁気記憶装置、半導体記憶装置、又は上記のいずれかの適切な組み合わせとすることができる。コンピュータ可読ストレージ媒体のより具体的な例の非網羅的なリストとして、以下のもの：即ち、ポータブル・コンピュータ・ディスク、ハード・ディスク、ランダム・アクセス・メモリ（RAM）、読み取り専用メモリ（ROM）、消去可能プログラム可能読み取り専用メモリ（EPROM又はフラッシュ・メモリ）、スタティック・ランダム・アクセス・メモリ（SRAM）、ポータブル・コンパクト・ディスク読み取り専用メモリ（CD-ROM）、デジタル多用途ディスク（DVD）、メモリ・スティック、フロッピー・ディスク、パンチカード又は命令が記録された溝内の隆起構造のような機械的にエンコードされたデバイス、及び上記のいずれかの適切な組み合わせが挙げられる。本明細書で使用される場合、コンピュータ可読ストレージ媒体は、電波、又は他の自由に伝搬する電磁波、導波管若しくは他の伝送媒体を通じて伝搬する電磁波（例えば、光ファイバ・ケーブルを通る光パルス）、又はワイヤを通して送られる電気信号などの、一時的信号自体として解釈されない。

【0040】

本明細書で説明されるコンピュータ可読プログラム命令は、コンピュータ可読ストレージ媒体からそれぞれのコンピューティング／処理デバイスに、又は、例えばインターネット、ローカル・エリア・ネットワーク、広域ネットワーク、及び／又は無線ネットワークなどのネットワークを介して外部コンピュータ又は外部ストレージ・デバイスにダウンロードすることができる。ネットワークは、銅伝送ケーブル、光伝送ファイバ、無線伝送、ルータ、ファイアウォール、スイッチ、ゲートウェイ・コンピュータ、及び／又はエッジ・サーバを含むことができる。各コンピューティング／処理デバイスにおけるネットワーク・アダプタ・カード又はネットワーク・インターフェースは、ネットワークからコンピュータ可読プログラム命令を受け取り、コンピュータ可読プログラム命令を転送して、それぞれのコンピューティング／処理デバイス内のコンピュータ可読ストレージ媒体内に格納する。

【0041】

本発明の動作を実行するためのコンピュータ可読プログラム命令は、アセンブラ命令、命令セット・アーキテクチャ（ISA）命令、マシン命令、マシン依存命令、マイクロコード、ファームウェア命令、状態設定データ、又は、Smalltalk、C++などのオブジェクト指向プログラミング言語、又は、「C」プログラミング言語若しくは類似のプログラミング言語などの通常の手続き型プログラミング言語を含む1つ又は複数のプログラミング言語の任意の組み合わせで記述することができる。コンピュータ可読プログラム命令は、完全にユーザのコンピュータ上で実行される場合もあり、一部がユーザのコンピュータ上で、独立型ソフトウェア・パッケージとして実行される場合もあり、一部がユーザのコンピュータ上で実行され、一部が遠隔コンピュータ上で実行される場合もあり、又は完全に遠隔コンピュータ若しくはサーバ上で実行される場合もある。最後のシナリオにおいては、遠隔コンピュータは、ローカル・エリア・ネットワーク（LAN）若しくは広域ネットワーク（WAN）を含むいずれかのタイプのネットワークを通じてユーザのコンピュータに接続される場合もあり、又は外部コンピュータへの接続がなされる場合もある（例えば、インターネット・サービス・プロバイダを用いたインターネットを通じて）。幾つかの実施形態において、例えば、プログラム可能論理回路、フィールド・プログラマブル・ゲート・アレイ（FPGA）、又はプログラム可能論理アレイ（PLA）を含む電子回路は、コンピュータ可読プログラム命令の状態情報を用いて、電子回路を個人化することによりコンピュータ可読プログラム命令を実行し、本発明の態様を実施することができる。

【0042】

本発明の態様は、本発明の実施形態による方法、装置（システム）及びコンピュータ・

10

20

30

40

50

プログラム製品のフローチャート図及び／又はブロック図を参照して説明される。フローチャート図及び／又はブロック図の各ブロック、並びにフローチャート図及び／又はブロック図内のブロックの組み合わせは、コンピュータ可読プログラム命令によって実装できることが理解されるであろう。

【0043】

これらのコンピュータ可読プログラム命令を、汎用コンピュータ、専用コンピュータ、又は他のプログラム可能データ処理装置のプロセッサに与えてマシンを製造し、それにより、コンピュータ又は他のプログラム可能データ処理装置のプロセッサによって実行される命令が、フローチャート及び／又はブロック図の1つ又は複数のブロック内で指定された機能／動作を実装するための手段を作り出すようにすることができる。これらのコンピュータ・プログラム命令を、コンピュータ、他のプログラム可能データ処理装置、又は他のデバイスを特定の方式で機能させるように指示することができるコンピュータ可読媒体内に格納し、それにより、そのコンピュータ可読媒体内に格納された命令が、フローチャート及び／又はブロック図の1つ又は複数のブロックにおいて指定された機能／動作を実装する命令を含む製品を製造するようにすることもできる。

10

【0044】

コンピュータ・プログラム命令を、コンピュータ、他のプログラム可能データ処理装置、又は他のデバイス上にロードして、一連の動作ステップをコンピュータ、他のプログラム可能データ処理装置、又は他のデバイス上で行わせてコンピュータ実施のプロセスを生成し、それにより、コンピュータ又は他のプログラム可能装置上で実行される命令が、フローチャート及び／又はブロック図の1つ又は複数のブロックにおいて指定された機能／動作を実行するためのプロセスを提供するようにすることもできる。

20

【0045】

図面内のフローチャート及びブロック図は、本発明の種々の実施形態による、システム、方法、及びコンピュータ・プログラム製品の可能な実装の、アーキテクチャ、機能及び動作を示す。この点に関して、フローチャート内の各ブロックは、指定された論理機能を実装するための1つ又は複数の実行可能命令を含む、モジュール、セグメント、又はコードの一部を表すことができる。幾つかの代替的な実装において、ブロック内に示される機能は、図に示される順序とは異なる順序で生じることがある。例えば、連続して示される2つのブロックは、関与する機能に応じて、実際には実質的に同時に実行されることもあり、又はこれらのブロックはときとして逆順で実行されることもある。ブロック図及び／又はフローチャート図の各ブロック、及びブロック図及び／又はフローチャート図内のブロックの組み合わせは、指定された機能又は動作を実行する、又は専用のハードウェアとコンピュータ命令との組み合わせを実行する、専用ハードウェア・ベースのシステムによって実装できることに留意されたい。

30

【符号の説明】

【0046】

- 10：アーキテクチャ
- 12、28：中央処理ユニット（CPU）
- 14：ストレージ・デバイス
- 16、26：メモリ・デバイス
- 18：通信ポート
- 20：通信ネットワーク
- 22：コンピュータ・システム
- 24：サーバ
- 30：大容量記憶装置
- 32：直接データ経路
- 34：トランスポート層
- 36：アプリケーション層
- 200：変換

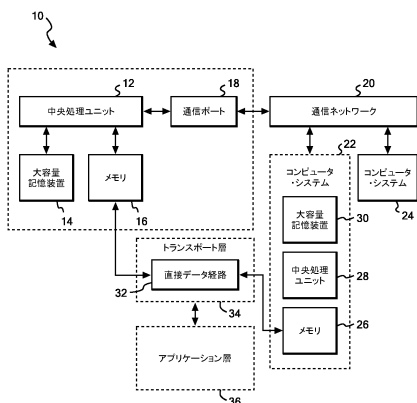
40

50

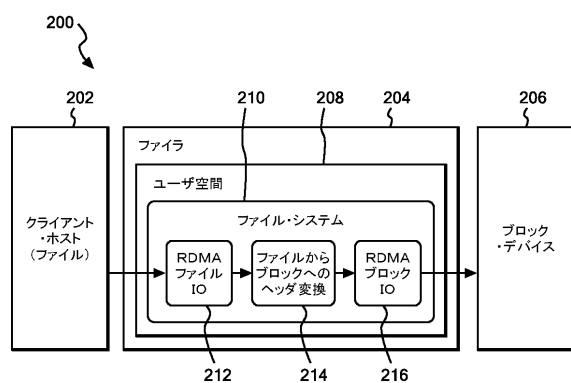
202 : ファイル・プロトコル・デバイス
 204 : ファイラ
 206、316 : ブロック・プロトコル・デバイス
 208、302、402 : ユーザ空間
 210 : ファイル・システム
 212 : ファイル入力 / 出力 (I / O) 機能
 214 : ファイルからブロックへのヘッダ変換操作
 216 : RDMA ブロック I / O 機能
 300、400 : 要求
 304、404 : カーネル空間
 306 : ファイルシステム・コア
 308、322 : リモート・データ・メモリ・アクセス (RDMA) 操作
 310、410 : ネットワーク・ファイル・システム (NFS) サーバ
 312 : 仮想ファイル・システム (VFS)
 314 : カーネル拡張
 408 : NFS RDMA 操作
 500、600 : 方法

10

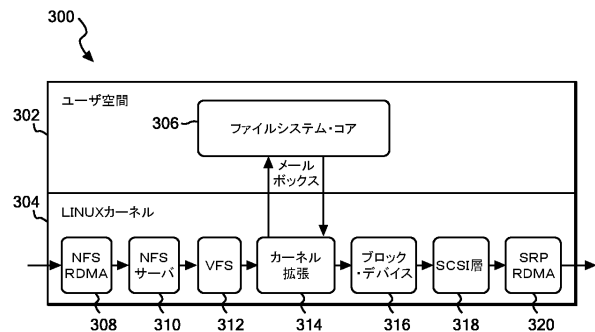
【図 1】



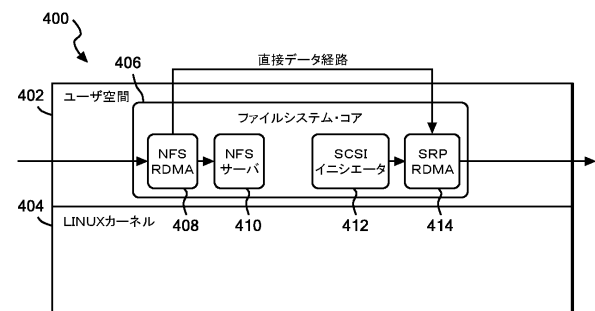
【図 2】



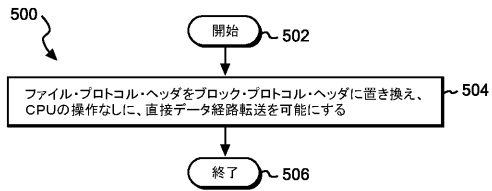
【図 3】



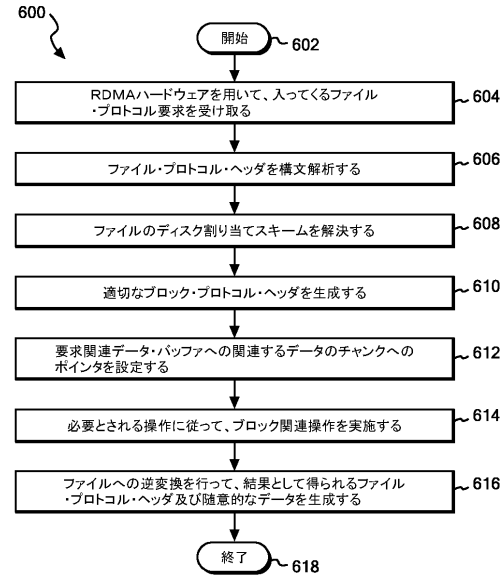
【図 4】



【図 5】



【図 6】



フロントページの続き

(74)代理人 100112690

弁理士 太佐 種一

(72)発明者 ドルーカー、ヴラディスラフ

アメリカ合衆国 8 5 7 4 4 アリゾナ州 ツーソン サウス・リタ・ロード 9 0 0 0

(72)発明者 アミット、ジョナサン

イスラエル国 6 7 2 0 1 テル・アビブ メナヘム・ベギン・ロード 1 3 2

(72)発明者 ロン、ザール

アメリカ合衆国 8 5 7 4 4 アリゾナ州 ツーソン サウス・リタ・ロード 9 0 0 0

(72)発明者 ローセン、ガル

イスラエル国 6 7 0 2 0 1 テル・アビブ メナヘム・ベギン・ロード 1 3 2

審査官 吉田 歩

(56)参考文献 特開2001-184294(JP,A)

特表2014-531685(JP,A)

米国特許出願公開第2012/0117341(US,A1)

特開2008-148181(JP,A)

(58)調査した分野(Int.Cl., DB名)

G 0 6 F 1 3 / 2 8

G 0 6 F 3 / 0 6