

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7160442号

(P7160442)

(45)発行日 令和4年10月25日(2022.10.25)

(24)登録日 令和4年10月17日(2022.10.17)

(51)国際特許分類

F I

G 0 6 F 16/188 (2019.01)

G 0 6 F 16/188

G 0 6 F 9/50 (2006.01)

G 0 6 F 9/50

1 2 0 Z

請求項の数 9 (全22頁)

(21)出願番号	特願2019-571200(P2019-571200)	(73)特許権者	390009531
(86)(22)出願日	平成30年6月14日(2018.6.14)		インターナショナル・ビジネス・マシー
(65)公表番号	特表2020-525909(P2020-525909		ンズ・コーポレーション
	A)		INTERNATIONAL BUSI
(43)公表日	令和2年8月27日(2020.8.27)		NESS MACHINES CORPO
(86)国際出願番号	PCT/IB2018/054378		RATION
(87)国際公開番号	WO2019/003029		アメリカ合衆国10504 ニューヨー
(87)国際公開日	平成31年1月3日(2019.1.3)		ク州 アーモンク ニュー オーチャード
審査請求日	令和2年11月30日(2020.11.30)		ロード
(31)優先権主張番号	15/636,770		New Orchard Road, A
(32)優先日	平成29年6月29日(2017.6.29)		rmonk, New York 105
(33)優先権主張国・地域又は機関			04, United States of
	米国(US)		America
(31)優先権主張番号	15/824,356	(74)代理人	100112690
(32)優先日	平成29年11月28日(2017.11.28)		弁理士 太佐 種一
	最終頁に続く		最終頁に続く

(54)【発明の名称】 読み取り / 書き込み要求を管理するための方法、コンピュータ・プログラムおよびコンピュータ・システム

(57)【特許請求の範囲】

【請求項1】

読み取り / 書き込み要求を管理するための方法であって、

テナント識別子のセット内の第1のテナント識別子に対応する第1のディレクトリを決定することであって、前記第1のディレクトリが第1の分散ファイル・システムを使用して構造化され、前記第1のテナント識別子が第1のテナントに対応する、前記決定することと、

接続サーバであるコネクタ・サービスを前記第1のテナントに対応する前記第1のディレクトリに割り当てることと、

前記コネクタ・サービスに対応する第2のディレクトリを決定することとであって、前記第2のディレクトリが、前記第1の分散ファイル・システムとは異なる第2の分散ファイル・システムを使用して構造化され、前記第2のディレクトリ上には、前記コネクタ・サービスとは別個に構成される第1のノードが含まれている第1のファイルのセットがあり、前記第1のファイルのセットが前記第1のテナントに対応するものである、前記決定することと、

前記コネクタ・サービスおよび前記第1のノードを使用して読み取り / 書き込み要求のセット内の第1の読み取り / 書き込み要求を処理することであって、前記コネクタ・サービスが、前記第1のテナントから来る、前記第1の分散ファイル・システムへの前記第1の読み取り / 書き込み要求を、前記第2の分散ファイル・システムに向けて方向付ける、前記処理することと、

10

20

前記第 1 の読み取り / 書き込み要求に対して第 1 の結果を生成することと
を含み、少なくとも前記コネクタ・サービスおよび前記第 1 のノードを使用して前記第 1 の読み取り / 書き込み要求を処理することが、コンピュータ・ハードウェア上で実行されるコンピュータ・ソフトウェアによって実行される、方法。

【請求項 2】

前記コネクタ・サービスと前記第 1 のノードに含まれる複数のテナントのそれぞれとを対応付けた経路を生成することをさらに含む、請求項 1 に記載の方法。

【請求項 3】

前記テナント識別子のセット内の第 2 のテナント識別子に対応する第 3 のディレクトリを決定することであって、前記第 3 のディレクトリが、前記第 1 の分散ファイル・システムを使用して構造化される、前記決定することと、

10

前記コネクタ・サービスを前記第 2 のテナント識別子に対応する前記第 3 のディレクトリに割り当てることと、

前記コネクタ・サービスおよび第 2 のノードを使用して第 2 の読み取り / 書き込み要求を処理することであって、前記第 2 のディレクトリ上には、前記第 2 のノードが含む第 2 のファイルのセットがあり、前記第 2 のファイルのセットが第 2 のテナントに対応するものである、前記処理することと、

前記第 2 の読み取り / 書き込み要求に対して第 2 の結果を生成することと
をさらに含んでいる、請求項 1 または 2 に記載の方法。

【請求項 4】

20

前記第 2 のノードが前記第 1 のノードである、請求項 3 に記載の方法。

【請求項 5】

前記第 1 の結果が、

新しいデータ・エントリと、

データのセットを含むメッセージと

から成る群から選択される、請求項 1 ないし 4 のいずれかに記載の方法。

【請求項 6】

前記第 2 の分散ファイル・システムは、P O S I X 互換である、請求項 1 ないし 5 のいずれかに記載の方法。

【請求項 7】

30

前記第 1 の分散ファイル・システムは、A p a c h e H a d o o p D i s t r i b u t e d F i l e S y s t e m (「HDFS」)である、請求項 1 ないし 6 のいずれかに記載の方法。

【請求項 8】

コンピュータ可読媒体に格納された、デジタル・コンピュータの内部メモリに読み込み可能なコンピュータ・プログラムであって、前記プログラムがコンピュータ上で実行された場合に請求項 1 ないし 7 のいずれかに記載の方法を実行するためのソフトウェア・コード部分を含んでいる、コンピュータ・プログラム。

【請求項 9】

読み取り / 書き込み要求を管理するためのコンピュータ・システムであって、

40

プロセッサ・セットと、

コンピュータ可読記憶媒体と

を備えており、

前記プロセッサ・セットが、前記コンピュータ可読記憶媒体に格納された命令を実行するように構造化されるか、配置されるか、接続されるか、またはプログラムされるか、あるいはその組み合わせが実行され、

前記命令が、

デバイスに、テナント識別子のセット内の第 1 のテナント識別子に対応する第 1 のディレクトリを決定させるように、前記デバイスによって実行可能なプログラム命令であって、前記第 1 のディレクトリが第 1 の分散ファイル・システムを使用して構造化され、前記

50

第 1 のテナント識別子が第 1 のテナントに対応する、前記プログラム命令と、

デバイスに、接続サーバであるコネクタ・サービスを前記第 1 のテナントに対応する前記第 1 のディレクトリに割り当てさせるように、前記デバイスによって実行可能なプログラム命令と、

デバイスに、前記コネクタ・サービスに対応する第 2 のディレクトリを決定させるように、前記デバイスによって実行可能なプログラム命令であって、前記第 2 のディレクトリが、前記第 1 の分散ファイル・システムとは異なる第 2 の分散ファイル・システムを使用して構造化され、前記第 2 のディレクトリ上には、前記コネクタ・サービスとは別個に構成される第 1 のノードが含んでいる第 1 のファイルのセットがあり、前記第 1 のファイルのセットが前記第 1 のテナントに対応するものである、前記プログラム命令と、

10

デバイスに、前記コネクタ・サービスおよび前記第 1 のノードを使用して読み取り / 書き込み要求のセット内の第 1 の読み取り / 書き込み要求を処理させるように、前記デバイスによって実行可能なプログラム命令であって、前記コネクタ・サービスが、前記第 1 のテナントから来る、前記第 1 の分散ファイル・システムへの前記第 1 の読み取り / 書き込み要求を、前記第 2 の分散ファイル・システムに向けて方向付ける、前記プログラム命令と、

デバイスに、前記第 1 の読み取り / 書き込み要求に対する第 1 の結果を生成させるように、前記デバイスによって実行可能なプログラム命令とを含んでいる、コンピュータ・システム。

【発明の詳細な説明】

20

【技術分野】

【0001】

本発明は、一般に、ストレージのアクセスおよび制御の分野に関し、特に、メモリの構成に関する。

【背景技術】

【0002】

集中型システムにおいて、仮想化は、計算リソース、ストレージ容量、またはアプリケーションのモビリティ、あるいはその組み合わせの弾力性を実現する。集中型インフラストラクチャは、情報技術のコンポーネントをソフトウェア・パッケージにグループ化する。仮想化コンテナは、ソフトウェアを信頼できる方法でサーバにインストールするためのファイル・システムを含むソフトウェア・パッケージである。仮想化コンテナの例は、Docker である。一部の仮想化コンテナは、ソフトウェア・ライブラリ・フレームワーク (software library frameworks) を含む。ソフトウェア・ライブラリ・フレームワークは、プログラミング・モデルを使用して、大規模なデータ・セットの分散処理を可能にする。そのようなソフトウェア・ライブラリ・フレームワークの一例は、Hadoop である。ポータブル・オペレーティング・システム・インターフェイスは、さまざまなオペレーティング・システム間の互換性を維持する。ポータブル・オペレーティング・システム・インターフェイスは、一連のアプリケーション・プログラミング・インターフェイスを定義する。ポータブル・オペレーティング・システム・インターフェイス規格の一例は、POSIX である。

30

【0003】

ビッグ・データ解析は、構造化データおよび非構造化データの両方を含むデータの指数関数的増加および可用性にもかかわらず、解析技術を可能にする。ビッグ・データ解析は、(i) 関係データベースに基づく超並列処理、および(ii) ソフトウェア・ライブラリ・フレームワークに基づく解析という 2 つの方向に発達した。

【0004】

さまざまなユーザの複数のクラスタを管理することは、極めて困難である。1 つの独立したネットワーク IP アドレスを監視しながら、1 つのテナントのクラスタ・インスタンスごとに 1 つのコネクタ・サービスを開始し、監視し、維持することができるが、この方法は、多量のシステム・リソースが必要になるため、拡張可能ではない。

40

50

【発明の概要】**【発明が解決しようとする課題】****【0005】**

したがって、当技術分野において、前述の問題に対処する必要がある。

【課題を解決するための手段】**【0006】**

第1の態様から見ると、本発明は、読み取り／書き込み要求を管理するための方法を提供し、この方法は、テナント識別子のセット内の第1のテナント識別子に対応する第1のディレクトリを決定することであって、第1のインターフェイス規格を使用して第1のディレクトリが構造化され、第1のテナント識別子が第1のディレクトリの第1のテナントに対応する、決定することと、コネクタ・サービスを第1のディレクトリおよび第1のテナント識別子に割り当てることと、コネクタ・サービスに対応する第2のディレクトリを決定することとであって、第2のインターフェイス規格を使用して第2のディレクトリが構造化され、第1のノードが第2のディレクトリ上に第1のファイルのセットを含んでおり、第1のファイルのセットが第1のテナントに対応する、決定することと、コネクタ・サービスおよび第1のノードを使用して読み取り／書き込み要求のセット内の第1の読み取り／書き込み要求を処理することであって、第1の読み取り／書き込み要求が第1のテナントから来る、処理することと、第1の読み取り／書き込み要求に対して第1の結果を生成することとであって、コネクタ・サービスおよび第1のノードを使用して少なくとも第1の読み取り／書き込み要求を処理することが、コンピュータ・ハードウェア上で実行されるコンピュータ・ソフトウェアによって実行される、生成することを含む。

【0007】

さらに別の態様から見ると、本発明は読み取り／書き込み要求を管理するためのコンピュータ・システムを提供し、このシステムは、プロセッサ・セットと、コンピュータ可読記憶媒体とを備えており、このプロセッサ・セットは、コンピュータ可読記憶媒体に格納された命令を実行するように構造化されるか、配置されるか、接続されるか、またはプログラムされるか、あるいはその組み合わせが実行され、それらの命令は、デバイスに、テナント識別子のセット内の第1のテナント識別子に対応する第1のディレクトリを決定させるように、デバイスによって実行可能なプログラム命令であって、第1のインターフェイス規格を使用して第1のディレクトリが構造化され、第1のテナント識別子が第1のディレクトリの第1のテナントに対応する、プログラム命令と、デバイスに、コネクタ・サービスを第1のディレクトリおよび第1のテナント識別子に割り当てさせるように、デバイスによって実行可能なプログラム命令と、デバイスに、コネクタ・サービスに対応する第2のディレクトリを決定させるように、デバイスによって実行可能なプログラム命令とであって、第2のインターフェイス規格を使用して第2のディレクトリが構造化され、第1のノードが第2のディレクトリ上に第1のファイルのセットを含んでおり、第1のファイルのセットが第1のテナントに対応する、プログラム命令と、デバイスに、コネクタ・サービスおよび第1のノードを使用して読み取り／書き込み要求のセット内の第1の読み取り／書き込み要求を処理させるように、デバイスによって実行可能なプログラム命令とであって、第1の読み取り／書き込み要求が第1のテナントから来る、プログラム命令と、デバイスに、第1の読み取り／書き込み要求に対する第1の結果を生成させるように、デバイスによって実行可能なプログラム命令とを含む。

【0008】

さらに別の態様から見ると、本発明は、読み取り／書き込み要求を管理するためのコンピュータ・プログラム製品を提供し、このコンピュータ・プログラム製品は、処理回路によって読み取り可能なコンピュータ可読記憶媒体を備えており、コンピュータ可読記憶媒体は、本発明のステップを実行するための方法を実行するためにこの処理回路によって実行される命令を格納している。

【0009】

さらに別の態様から見ると、本発明は、コンピュータ可読媒体に格納された、デジタル

・コンピュータの内部メモリに読み込み可能なコンピュータ・プログラムを提供し、このコンピュータ・プログラムは、コンピュータ上で実行された場合に本発明のステップを実行するためのソフトウェア・コード部分を含んでいる。

【 0 0 1 0 】

本発明の態様によれば、方法、コンピュータ・プログラム製品、またはシステム、あるいはその組み合わせが存在し、これらは、(i) テナント識別子のセット内の第 1 のテナント識別子に対応する第 1 のディレクトリを決定することであって、(a) 第 1 のインターフェイス規格を使用して第 1 のディレクトリが構造化され、(b) 第 1 のテナント識別子が第 1 のディレクトリの第 1 のテナントに対応する、決定することと、(i i) コネクタ・サービスを第 1 のディレクトリおよび第 1 のテナント識別子に割り当てることと、(i i i) コネクタ・サービスに対応する第 2 のディレクトリを決定することであって、(a) 第 2 のインターフェイス規格を使用して第 2 のディレクトリが構造化され、(b) 第 1 のノードが第 2 のディレクトリ上に第 1 のファイルのセットを含んでおり、(c) 第 1 のファイルのセットが第 1 のテナントに対応する、決定することと、(i v) コネクタ・サービスおよび第 1 ノードを使用して読み取り / 書き込み要求のセット内の第 1 の読み取り / 書き込み要求を処理することであって、第 1 の読み取り / 書き込み要求が第 1 のテナントから来る、処理することと、(v) 第 1 の読み取り / 書き込み要求に対して第 1 の結果を生成することとを実行する動作を(必ずしもこの順序でなく)実行する。コネクタ・サービスおよび第 1 のノードを使用して少なくとも第 1 の読み取り / 書き込み要求を処理することは、コンピュータ・ハードウェア上で実行されるコンピュータ・ソフトウェアによって実行される。

【 0 0 1 1 】

以下では、次の図に示された好ましい実施形態を単に例として参照し、本発明が説明される。

【図面の簡単な説明】

【 0 0 1 2 】

【図 1】本発明に従うシステムの第 1 の実施形態を示すブロック図である。

【図 2】第 1 の実施形態のシステムによって少なくとも一部が実行される第 1 の実施形態の方法を示すフローチャートである。

【図 3】第 1 の実施形態のシステムの機械論理(例えば、ソフトウェア)部分を示すブロック図である。

【図 4】本発明に従うシステムの第 2 の実施形態によって実行される第 2 の実施形態の方法を示すフローチャートである。

【図 5】システムの第 2 の実施形態を示すブロック図である。

【図 6】本発明に従うシステムの第 3 の実施形態によって生成されたルックアップ・テーブルである。

【図 7】本発明に従うシステムの第 4 の実施形態によって実行される第 3 の実施形態の方法を示すフローチャートである。

【発明を実施するための形態】

【 0 0 1 3 】

ノード上のマルチテナント分散ファイル・システムの構成。さまざまなテナントおよびテナント・クラスタ(tenant clusters)は、分散ファイル・システムと相互関係があり、分散ファイル・システムは、コネクタ・サービスを介してさまざまなテナントと通信する。分散ファイル・システム全体は、物理ノード上に存在する。「発明を実施するための形態」のセクションは、(i) ハードウェアおよびソフトウェア環境、(i i) 実施形態例、(i i i) 追加のコメントまたは実施形態あるいはその両方、ならびに(i v) 定義のサブセクションに分割されている。

【 0 0 1 4 】

I . ハードウェアおよびソフトウェア環境

本発明は、任意の可能な統合の技術的詳細レベルで、システム、方法、またはコンピュ

10

20

30

40

50

ータ・プログラム製品、あるいはその組み合わせであってよい。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行させるためのコンピュータ可読プログラム命令を含んでいるコンピュータ可読記憶媒体を含んでよい。

【0015】

コンピュータ可読記憶媒体は、命令実行デバイスによって使用するための命令を保持および格納できる有形のデバイスにすることができる。コンピュータ可読記憶媒体は、例えば、電子ストレージ・デバイス、磁気ストレージ・デバイス、光ストレージ・デバイス、電磁ストレージ・デバイス、半導体ストレージ・デバイス、またはこれらの任意の適切な組み合わせであってよいが、これらに限定されない。コンピュータ可読記憶媒体のさらに具体的な例の非網羅的リストは、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: random access memory)、読み取り専用メモリ (ROM: read-only memory)、消去可能プログラマブル読み取り専用メモリ (EPROM: erasable programmable read-only memory またはフラッシュ・メモリ)、スタティック・ランダム・アクセス・メモリ (SRAM: static random access memory)、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: compact disc read-only memory)、デジタル多用途ディスク (DVD: digital versatile disk)、メモリ・スティック、フロッピー (R)・ディスク、パンチカードまたは命令が記録されている溝の中の隆起構造などの機械的にエンコードされるデバイス、およびこれらの任意の適切な組み合わせを含む。本明細書において使用されるとき、コンピュータ可読記憶媒体は、それ自体が、電波またはその他の自由に伝搬する電磁波、導波管またはその他の送信媒体を伝搬する電磁波 (例えば、光ファイバ・ケーブルを通過する光パルス)、あるいはワイヤを介して送信される電気信号などの一過性の信号であると解釈されるべきではない。

【0016】

本明細書に記載されたコンピュータ可読プログラム命令は、コンピュータ可読記憶媒体から各コンピューティング・デバイス / 処理デバイスへ、またはネットワーク (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、または無線ネットワーク、あるいはその組み合わせ) を介して外部コンピュータまたは外部ストレージ・デバイスへダウンロードされ得る。このネットワークは、銅伝送ケーブル、光伝送ファイバ、無線送信、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、またはエッジ・サーバ、あるいはその組み合わせを備えてよい。各コンピューティング・デバイス / 処理デバイス内のネットワーク・アダプタ・カードまたはネットワーク・インターフェイスは、コンピュータ可読プログラム命令をネットワークから受信し、それらのコンピュータ可読プログラム命令を各コンピューティング・デバイス / 処理デバイス内のコンピュータ可読記憶媒体に格納するために転送する。

【0017】

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: instruction-set-architecture) 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、集積回路のための構成データ、あるいは、Smalltalk (R)、C++ などのオブジェクト指向プログラミング言語、および「C」プログラミング言語または同様のプログラミング言語などの手続き型プログラミング言語を含む1つまた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されたソース・コードまたはオブジェクト・コードであってよい。コンピュータ可読プログラム命令は、ユーザのコンピュータ上で全体的に実行すること、ユーザのコンピュータ上でスタンドアロン・ソフトウェア・パッケージとして部分的に実行すること、ユーザのコンピュータ上およびリモート・コンピュータ上でそれぞれ部分的に実行すること、あるいはリモート・コンピュータ上またはサーバ上で全体的に実行することができる。後者のシナリオでは、リモート・コンピュータを、ローカル・エリア・ネットワーク (LAN: local area network) または広域ネットワーク (WAN: wide areanetwork) を含む任意の種類のネットワークを介してユーザのコンピュータに接続する

10

20

30

40

50

ことができ、または接続を、（例えば、インターネット・サービス・プロバイダを使用してインターネットを介して）外部コンピュータに対して行うことができる。一部の実施形態では、本発明の態様を実行するために、例えばプログラマブル論理回路、フィールドプログラマブル・ゲート・アレイ（FPGA：field-programmable gate arrays）、またはプログラマブル・ロジック・アレイ（PLA：programmable logic arrays）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利用することによって、電子回路をカスタマイズするためのコンピュータ可読プログラム命令を実行してよい。

【0018】

本発明の態様は、本明細書において、本発明の実施形態に従って、方法、装置（システム）、およびコンピュータ・プログラム製品のフローチャート図またはブロック図あるいはその両方を参照して説明される。フローチャート図またはブロック図あるいはその両方の各ブロック、ならびにフローチャート図またはブロック図あるいはその両方に含まれるブロックの組み合わせが、コンピュータ可読プログラム命令によって実装され得るということが理解されるであろう。

【0019】

これらのコンピュータ可読プログラム命令は、コンピュータまたはその他のプログラム可能なデータ処理装置のプロセッサを介して実行される命令が、フローチャートまたはブロック図あるいはその両方の1つまたは複数のブロックに指定される機能／動作を実施する手段を作り出すべく、汎用コンピュータ、専用コンピュータ、または他のプログラム可能なデータ処理装置のプロセッサに提供されてマシンを作り出すものであってよい。これらのコンピュータ可読プログラム命令は、命令が格納されたコンピュータ可読記憶媒体がフローチャートまたはブロック図あるいはその両方の1つまたは複数のブロックに指定される機能／動作の態様を実施する命令を含んでいる製品を備えるように、コンピュータ可読記憶媒体に格納され、コンピュータ、プログラム可能なデータ処理装置、または他のデバイス、あるいはその組み合わせに特定の方式で機能するように指示できるものであってよい。

【0020】

コンピュータ可読プログラム命令は、コンピュータ上、その他のプログラム可能な装置上、またはその他のデバイス上で実行される命令が、フローチャートまたはブロック図あるいはその両方の1つまたは複数のブロックに指定される機能／動作を実施するように、コンピュータ実装プロセスを生成すべく、コンピュータ、その他のプログラム可能なデータ処理装置、またはその他のデバイスに読み込まれて、コンピュータ上、その他のプログラム可能な装置上、またはその他のデバイス上で一連の動作可能なステップを実行させるものであってよい。

【0021】

図内のフローチャートおよびブロック図は、本発明のさまざまな実施形態に従って、システム、方法、およびコンピュータ・プログラム製品の可能な実装のアーキテクチャ、機能、および動作を示す。これに関連して、フローチャートまたはブロック図内の各ブロックは、規定された論理機能を実装するための1つまたは複数の実行可能な命令を備える、命令のモジュール、セグメント、または部分を表してよい。一部の代替の実装では、ブロックに示された機能は、図に示された順序とは異なる順序で発生してよい。例えば、連続して示された2つのブロックは、実際には、含まれている機能に応じて、実質的に同時に実行されるか、または場合によっては逆の順序で実行されてよい。ブロック図またはフローチャート図あるいはその両方の各ブロック、ならびにブロック図またはフローチャート図あるいはその両方に含まれるブロックの組み合わせは、規定された機能または動作を実行するか、または専用ハードウェアとコンピュータ命令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され得るということにも注意する。

【0022】

ここで、本発明に従うソフトウェアまたは方法あるいはその両方の可能なハードウェアおよびソフトウェア環境の実施形態が、各図を参照して詳細に説明される。図1は、マル

10

20

30

40

50

チテナント構成サブシステム 102、ユーザ・サブシステム 104、仮想コンテナ・サブシステム 106、仮想コンテナ・サブシステム 108、コネクタ・サービス 112、および通信ネットワーク 114を含んでいるネットワーク・コンピュータ・システム 100のさまざまな部分を示す機能ブロック図である。マルチテナント構成サブシステム 102は、マルチテナント構成コンピュータ 200、ディスプレイ・デバイス 212、および外部デバイス 214を含んでいる。マルチテナント構成コンピュータ 200は、通信ユニット 202、プロセッサ・セット 204、入出力 (I/O: input/output) インターフェイス・セット 206、メモリ・デバイス 208、および永続的ストレージ・デバイス 210を含んでいる。メモリ・デバイス 208は、ランダム・アクセス・メモリ (RAM) デバイス 216およびキャッシュ・メモリ・デバイス 218を含んでいる。永続的ストレージ 210は、マルチテナント構成プログラム 300を含んでいる。仮想コンテナ・サブシステム 108は、ソフトウェア・ライブラリ・フレームワーク 110を含んでいる。

10

【0023】

マルチテナント構成サブシステム 102は、多くの点で、本発明におけるさまざまなコンピュータ・サブシステムを代表している。したがって以下では、マルチテナント構成サブシステム 102の複数の部分が説明される。

【0024】

マルチテナント構成サブシステム 102は、ラップトップ・コンピュータ、タブレット・コンピュータ、ネットブック・コンピュータ、パーソナル・コンピュータ (PC: personal computer)、デスクトップ・コンピュータ、パーソナル・デジタル・アシスタント (PDA: personal digital assistant)、スマートフォン、または通信ネットワーク 114を介してクライアント・サブシステムと通信できる任意のプログラム可能な電子デバイスであってよい。マルチテナント構成プログラム 300は、以下で、この「発明を実施するための形態」セクションの「実施形態例」サブセクションにおいて詳細に説明される、特定のソフトウェア機能を作成し、管理し、制御するために使用される機械可読の命令またはデータあるいはその両方の集合である。

20

【0025】

マルチテナント構成サブシステム 102は、通信ネットワーク 114を介して他のコンピュータ・サブシステムと通信することができる。例えば、通信ネットワーク 114は、ローカル・エリア・ネットワーク (LAN)、インターネットなどの広域ネットワーク (WAN)、またはこれらの組み合わせであることができ、有線接続、無線接続、または光ファイバ接続を含むことができる。一般に、通信ネットワーク 114は、マルチテナント構成サブシステム 102とクライアント・サブシステムの間の通信をサポートする接続およびプロトコルの任意の組み合わせであることができる。

30

【0026】

マルチテナント構成サブシステム 102は、多くの両方向の矢印を含むブロック図として示されている。これらの両方向の矢印 (別の参照番号はない) は、マルチテナント構成サブシステム 102のさまざまなコンポーネント間の通信を提供する通信ファブリックを表す。通信ファブリックは、プロセッサ (マイクロプロセッサ、通信プロセッサ、またはネットワーク・プロセッサ、あるいはその組み合わせなど)、システム・メモリ、周辺機器、およびシステム内の任意のその他のハードウェア・コンポーネントの間で、データまたは制御情報あるいはその両方を渡すために設計された、任意のアーキテクチャを使用して実装され得る。例えば、通信ファブリックは、少なくとも一部において、1つまたは複数のバスを使用して実装され得る。

40

【0027】

メモリ・デバイス 208および永続的ストレージ・デバイス 210は、コンピュータ可読記憶媒体である。一般に、メモリ・デバイス 208は、任意の適切な揮発性または不揮発性のコンピュータ可読記憶媒体を含むことができる。現在または近い将来あるいはその両方において、(i) 外部デバイス 214が、マルチテナント構成サブシステム 102の一部または全部のメモリを提供できてよいということ、または (ii) マルチテナント構

50

成サブシステム 102 の外部にあるデバイスが、マルチテナント構成サブシステム 102 にメモリを提供できてよいということ、あるいはその両方のことに、さらに注意する。

【0028】

マルチテナント構成プログラム 300 は、プロセッサ・セット 204 の 1 つまたは複数のプロセッサによる、通常はメモリ・デバイス 208 を介したアクセスまたは実行あるいはその両方のために、永続的ストレージ・デバイス 210 に格納される。永続的ストレージ・デバイス 210 は、(i) 少なくとも送信中の信号より永続的であり、(ii) プログラム (ソフト・ロジックまたはデータあるいはその両方を含む) を有形の媒体 (磁気または光領域など) に格納し、(iii) 永久記憶装置より大幅に永続性が低い。代替としてデータ格納は、永続的ストレージ・デバイス 210 によって提供される格納のタイプより永続的または永久的あるいはその両方であってよい。

10

【0029】

マルチテナント構成プログラム 300 は、実質的データ (すなわち、データベースに格納されるデータのタイプ) または機械可読の実行可能な命令あるいはその両方を含んでよい。この特定の実施形態 (すなわち、図 1) では、永続的ストレージ・デバイス 210 は、磁気ハード・ディスク・ドライブを含んでいる。可能性のある変形をいくつか挙げると、永続的ストレージ・デバイス 210 は、固体ハード・ドライブ、半導体ストレージ・デバイス、読み取り専用メモリ (ROM)、消去可能プログラマブル読み取り専用メモリ (EPROM)、フラッシュ・メモリ、あるいはプログラム命令またはデジタル情報を格納できる任意のその他のコンピュータ可読記憶媒体を含んでよい。

20

【0030】

永続的ストレージ・デバイス 210 によって使用される媒体は、取り外し可能であってよい。例えば、取り外し可能ハード・ドライブを、永続的ストレージ・デバイス 210 に使用できる。その他の例としては、永続的ストレージ・デバイス 210 の一部でもある別のコンピュータ可読記憶媒体に転送するためのドライブに挿入される、光ディスクおよび磁気ディスク、サム・ドライブ、ならびにスマート・カードが挙げられる。

【0031】

これらの例において、通信ユニット 202 は、マルチテナント構成サブシステム 102 の外部にある他のデータ処理システムまたはデバイスとの通信を提供する。これらの例において、通信ユニット 202 は、1 つまたは複数のネットワーク・インターフェイス・カードを含む。通信ユニット 202 は、物理的通信リンクまたは無線通信リンクのいずれか、あるいはその両方を使用して通信を提供できる。本明細書において説明される任意のソフトウェア・モジュールは、通信ユニット (通信ユニット 202 など) を介して永続的ストレージ・デバイス (永続的ストレージ・デバイス 210 など) にダウンロードされてよい。

30

【0032】

I/O インターフェイス 206 は、データ通信でマルチテナント構成コンピュータ 200 にローカルに接続されることがある他のデバイスとのデータの入力および出力を可能にする。例えば、I/O インターフェイス 206 は、外部デバイス 214 との接続を提供する。外部デバイス 214 は通常、キーボード、キーパッド、タッチ・スクリーン、またはその他の適切な入力デバイス、あるいはその組み合わせなどのデバイスを含む。外部デバイス 214 は、例えばサム・ドライブ、ポータブル光ディスクまたはポータブル磁気ディスク、およびメモリ・カードなどの、ポータブル・コンピュータ可読記憶媒体を含むこともできる。本発明の実施形態 (例えば、マルチテナント構成プログラム 300) を実践するために使用されるソフトウェアおよびデータは、そのようなポータブル・コンピュータ可読記憶媒体に格納され得る。それらの実施形態では、関連するソフトウェアが、I/O インターフェイス・セット 206 を介して、全体的または部分的に永続的ストレージ・デバイス 210 に読み込まれてよい (または、読み込まれなくてよい)。I/O インターフェイス・セット 206 は、データ通信でディスプレイ・デバイス 212 にも接続される。

40

【0033】

50

ディスプレイ・デバイス 212 は、データをユーザに表示するためのメカニズムを提供し、例えば、コンピュータのモニタまたはスマートフォンの表示画面であってよい。

【0034】

本明細書に記載されたプログラムは、アプリケーションに基づいて識別され、本発明の特定の実施形態において、そのアプリケーションに関して実装される。ただし、本明細書における特定のプログラムの名前は単に便宜上使用されていると理解されるべきであり、したがって、本発明は、そのような名前によって識別されたか、または暗示されたか、あるいはその両方によって示された特定のアプリケーションのみで使用するよう制限されるべきではない。

【0035】

本発明のさまざまな実施形態の説明は、例示の目的で提示されているが、網羅的であることは意図されておらず、開示された実施形態に制限されない。説明された実施形態の範囲を逸脱することなく多くの変更および変形が可能であることは、当業者にとって明らかである。本明細書で使用された用語は、実施形態の原理、実際の適用、または市場で見られる技術を超える技術的改良を最も適切に説明するため、または他の当業者が本明細書で開示された実施形態を理解できるようにするために選択されている。

【0036】

II. 実施形態例

図2は、本発明に従う方法を表すフローチャート250を示している。図3は、フローチャート250の方法の動作の少なくとも一部を実行するマルチテナント構成プログラム300を示している。以降では、図2（方法の動作のブロックに関する）および図3（ソフトウェア・ブロックに関する）を詳細に参照して、この方法および関連するソフトウェアについて説明する。

【0037】

動作S255で処理が開始し、要求受信モジュール（「mod」）302が要求のセットを受信する。本発明の一部の実施形態では、要求受信mod302が、要求元のセットから要求のセットを受信する。要求元の例としては、ソフトウェア・ライブラリ・フレームワーク、仮想コンテナ、またはユーザ、あるいはその組み合わせが挙げられるが、これらに限定されない。一部の実施形態では、要求のセットは、入出力（「I/O」）要求のセットである。さらに別の実施形態では、要求のセットは、読み取り/書き込み要求のセットである。これらの実施形態の一部では、要求のセットは、I/O読み取り/書き込み要求のセットである。仮想コンテナの例は、Dockerである。ソフトウェア・ライブラリ・フレームワークの例は、Hadoopである。さらに別の実施形態では、要求受信mod302が、要求元の動的インスタンス化のセットから要求のセットを受信する。

【0038】

一部の実施形態では、要求元は第1の分散ファイル・システムである。これらの実施形態の一部では、第1の分散ファイル・システムはPOSIX互換である。さらに別の実施形態では、第1の分散ファイル・システムは、第1のインターフェイス規格を使用して構造化される。一部の実施形態では、要求のセットは第2の分散ファイル・システムに関連している。これらの実施形態の一部では、第2の分散ファイル・システムはPOSIX互換ではない。さらに別の実施形態では、第2の分散ファイル・システムは、第2のインターフェイス規格を使用して構造化される。代替として、一部の実施形態では、（i）第1の分散ファイル・システムがPOSIX互換であり、（ii）第2の分散ファイル・システムがPOSIX互換ではない。さらに別の代替の実施形態では、第1の分散ファイル・システムも第2の分散ファイル・システムもPOSIX互換ではなく、異なるインターフェイス規格を使用して第1の分散ファイル・システムおよび第2の分散ファイル・システムが構造化される。

【0039】

処理が動作S260に進み、ディレクトリ決定mod304が、要求元のセットに対応するディレクトリのセットを決定する。本発明の一部の実施形態では、ディレクトリ決定

10

20

30

40

50

mod 304 が、要求元のセットに対応するディレクトリのセットを決定する。ディレクトリは、コンピュータ・ファイルのセットを整理するための構造である。ディレクトリは、パス、フォルダ、またはドローと呼ばれることもあり、あるいはその組み合わせで呼ばれることもある。ディレクトリは、(i) 親フォルダ / 子フォルダ / ファイル、拡張子、または (i i) 親フォルダ > 子フォルダ > ファイル、あるいはその両方を含む、さまざまな形態で表され得る。これらの実施形態の一部では、ディレクトリ決定 mod 304 が、テナント識別子のセットに対応するディレクトリのセットを決定する。他の実施形態では、ディレクトリ決定 mod 304 が、ディレクトリを要求元のセットに割り当てることによって、テナント識別子のセットに対応するディレクトリのセットを決定する。さらに別の実施形態では、ディレクトリ決定 mod 304 が、サブディレクトリを要求元のセットに割り当てることによって、テナント識別子のセットに対応するディレクトリのセットを決定する。一部の実施形態では、要求元のセット内の第 1 の要求元が、第 1 のディレクトリに対応する。他の実施形態では、要求元のセットが、第 1 のディレクトリを共有する。一部の実施形態では、ディレクトリ決定 mod 304 が、要求受信 mod 302 が動作 S 255 で受信した要求のセットの送信元である、要求元のセットに対応するディレクトリのセットを決定する。

【 0040 】

処理が動作 S 265 に進み、テナント識別子決定 mod 306 が、要求のセットに対応するテナント識別子のセットを決定する。本発明の一部の実施形態では、テナント識別子決定 mod 306 が、要求のセットに対応するテナント識別子のセットを決定する。一部の実施形態では、テナント識別子決定 mod 306 が、動的インスタンス化である要求元のセットのテナント識別子のセットを決定する。代替の実施形態では、テナント識別子決定 mod 306 が、仮想コンテナのセットのテナント識別子のセットを決定する。さらに別の実施形態では、テナント識別子決定 mod 306 が、ソフトウェア・ライブラリ・フレームワークのセットのテナント識別子のセットを決定する。代替として、テナント識別子決定 mod 306 が、ユーザのセットのテナント識別子のセットを決定する。一部の実施形態では、テナント識別子決定 mod 306 が、テナントのセットのインスタンスのセットのテナント識別子のセットを決定する。一部の実施形態では、テナント識別子決定 mod 306 が、動作 S 255 で要求受信 mod 302 によって受信された要求のセットに対応するテナント識別子のセットを決定する。代替として、テナント識別子決定 mod 306 が、動作 S 260 でディレクトリ決定 mod 304 によって決定されたディレクトリのセットに対応するテナント識別子のセットを決定する。

【 0041 】

処理が動作 S 270 に進み、コネクタ・サービス割り当て mod 308 がコネクタ・サービスを割り当てる。本発明の一部の実施形態では、コネクタ・サービス割り当て mod 308 がコネクタ・サービスを割り当てる。さらに別の実施形態では、コネクタ・サービスはコンピュータ・システム上の唯一のコネクタ・サービスである。代替として、コネクタ・サービスは、第 1 の分散ファイル・システムおよび第 2 の分散ファイル・システムに関連付けられた唯一のコネクタ・サービスである。これらの実施形態の一部では、コネクタ・サービスが、第 1 の分散ファイル・システム上の要求元のセットからの要求を、第 2 の分散ファイル・システムに向けて方向付ける。他の実施形態では、コネクタ・サービス割り当て mod 308 が、少なくとも一部において、テナント識別子のセットに基づいて、コネクタ・サービスを割り当てる。さらに別の実施形態では、コネクタ・サービス割り当て mod 308 が、少なくとも一部において、ディレクトリのセットに基づいて、コネクタ・サービスを割り当てる。コネクタ・サービスは、接続サーバと呼ばれることもある。コネクタ・サービスは、適切なチャネルのセットを介して、要求のセットを方向付ける。接続サーバは、(i) ユーザのセットを認証すること、(i i) リソースのセットに対する資格をユーザのセットに与えること、(i i i) パッケージのセットをリソースのセットに割り当てること、(i v) ローカル・セッションもしくはリモート・セッションまたはその両方を管理すること、(v) セキュリティで保護された接続のセットを確立する

10

20

30

40

50

こと、または(v i)ポリシーを適用すること、あるいはその組み合わせを実行することを含むが、これらに限定されない、機能を実行してもよい。一部の実施形態では、コネクタ・サービス割り当てmod 308が、少なくとも一部において、動作S 255で要求受信mod 302によって受信された要求のセットの要求元のセットに基づいて、コネクタ・サービスを割り当てる。他の実施形態では、コネクタ・サービス割り当てmod 308が、少なくとも一部において、動作S 255で要求受信mod 302によって受信された要求のセットに基づいて、コネクタ・サービスを割り当てる。さらに別の実施形態では、コネクタ・サービス割り当てmod 308が、少なくとも一部において、動作S 260でディレクトリ決定mod 304によって決定されたディレクトリのセットに基づいて、コネクタ・サービスを割り当てる。代替の実施形態では、コネクタ・サービス割り当てmod 308が、少なくとも一部において、動作S 265でテナント識別子決定mod 306によって決定されたテナント識別子のセットに基づいて、コネクタ・サービスを割り当てる。

10

【0042】

処理が動作S 275に進み、ノード決定mod 310が、要求元のセットに対応するノードを決定する。本発明の一部の実施形態では、ノード決定mod 310が、要求元のセットに対応するノードを決定する。これらの実施形態の一部では、ノード決定mod 310が、第1のノードが要求元のセット内の各要求元に対応するということを決定する。これらの実施形態の一部では、ノード決定mod 310が、物理ノードが要求元のセットに対応するということを決定する。他の実施形態では、ノード決定mod 310が、仮想ノードが要求元のセットに対応するということを決定する。代替の実施形態では、ノード決定mod 310が、要求元のセット内の各要求元を第1のノードに割り当てることによって、要求元のセットに対応するノードを決定する。一部の実施形態では、ノード決定mod 310が、要求のセットに対応するノードを決定する。さらに別の実施形態では、ノード決定mod 310が、テナント識別子のセットに対応するノードを決定する。他の実施形態では、ノード決定mod 310が、少なくとも一部において、コネクタ・サービスに基づいて、ノードを決定する。代替の実施形態では、ノード決定mod 310が、少なくとも一部において、ノードとコネクタ・サービスの間の1対1の関係に基づいて、ノードを決定する。他の実施形態では、ノード決定mod 310が、コネクタ・サービスとノードの間の経路をマッピングする。一部の実施形態では、ノード決定mod 310が、要求受信mod 302が動作S 255で受信した要求のセットの送信元である、要求元のセットに対応するノードを決定する。他の実施形態では、ノード決定mod 310が、動作S 255で要求受信mod 302によって受信された要求のセットに対応するノードを決定する。さらに別の実施形態では、ノード決定mod 310が、動作S 260でディレクトリ決定mod 304によって決定されたディレクトリのセットに対応するノードを決定する。代替の実施形態では、ノード決定mod 310が、動作S 265でテナント識別子決定mod 306によって決定されたテナント識別子のセットに対応するノードを決定する。代替として、ノード決定mod 310が、少なくとも一部において、動作S 270でコネクタ・サービス割り当てmod 308によって割り当てられたコネクタ・サービスに基づいて、ノードを決定する。

20

30

40

【0043】

処理が動作S 280に進み、要求処理mod 312が要求のセットを処理する。本発明の一部の実施形態では、要求処理mod 312が要求のセットを処理する。一部の実施形態では、要求処理mod 312が、少なくとも一部において、テナント識別子のセットに基づいて、要求のセットを処理する。他の実施形態では、要求処理mod 312が、少なくとも一部において、ノードに基づいて、要求のセットを処理する。さらに別の実施形態では、要求処理mod 312が、少なくとも一部において、ディレクトリに基づいて、要求のセットを処理する。一部の実施形態では、要求処理mod 312が、第1の分散ファイル・システムを第2の分散ファイル・システムにマウントする。代替の実施形態では、要求処理mod 312が、少なくとも一部において、コネクタ・サービスに基づいて、要

50

求のセットを処理する。読み取り要求の場合、要求処理 `mod 3 1 2` が、データのセットをストレージから読み取る。書き込み要求の場合、要求処理 `mod 3 1 2` が、ストレージ内のデータのセットを変更する。入力要求の場合、要求処理 `mod 3 1 2` が、データのセットを受信する。出力要求の場合、要求処理 `mod 3 1 2` が、データのセットを送信する。一部の実施形態では、要求処理 `mod 3 1 2` が、動作 `S 2 5 5` で要求受信 `mod 3 1 2` によって受信された要求のセットを処理する。他の実施形態では、要求処理 `mod 3 1 2` が、少なくとも一部において、動作 `S 2 6 5` でテナント識別子決定 `mod 3 0 6` によって決定されたテナント識別子のセットに基づいて、要求のセットを処理する。さらに別の実施形態では、要求処理 `mod 3 1 2` が、少なくとも一部において、動作 `S 2 7 5` でノード決定 `mod 3 1 0` によって決定されたノードに基づいて、要求のセットを処理する。他の実施形態では、要求処理 `mod 3 1 2` が、少なくとも一部において、動作 `S 2 6 0` でディレクトリ決定 `mod 3 0 4` によって決定されたディレクトリのセットに基づいて、要求のセットを処理する。代替の実施形態では、要求処理 `mod 3 1 2` が、少なくとも一部において、動作 `S 2 7 0` でコネクタ・サービス決定 `mod 3 0 8` によって決定されたコネクタ・サービスに基づいて、要求のセットを処理する。

【 0 0 4 4 】

動作 `S 2 8 5` で処理が終了し、結果生成 `mod 3 1 4` が結果のセットを生成する。本発明の一部の実施形態では、結果生成 `mod 3 1 4` が、要求のセットの結果のセットを生成する。一部の実施形態では、結果生成 `mod 3 1 4` が、データのセットを含んでいるメッセージのセットを生成することによって、読み取り要求のセットに対する結果のセットを生成する。一部の実施形態では、結果生成 `mod 3 1 4` が、新しいデータ・エントリのセットを生成することによって、書き込み要求のセットに対する結果のセットを生成する。一部の実施形態では、結果生成 `mod 3 1 4` が、受信されたデータのセットを格納することによって、入力要求のセットに対する結果のセットを生成する。一部の実施形態では、結果生成 `mod 3 1 4` が、メッセージのセットを生成することによって、出力要求のセットに対する結果のセットを生成する。他の実施形態では、結果生成 `mod 3 1 4` が、`POSI X` 互換ではない第 1 の分散ファイル・システムの結果を生成する。さらに別の実施形態では、結果生成 `mod 3 1 4` が、`Hadoop` である第 1 の分散ファイル・システムの結果のセットを生成する。他の実施形態では、結果は、新しいデータ・エントリまたはデータのセットを含むメッセージあるいはその両方を含むが、これらに限定されない。一部の実施形態では、結果生成 `mod 3 1 4` が、動作 `S 2 5 5` で要求受信 `mod 3 0 2` によって受信された要求のセットに対する結果のセットを生成する。

【 0 0 4 5 】

III. 追加のコメントまたは実施形態あるいはその両方

本発明の一部の実施形態は、現在の技術に対して改善を行うために、次の要因、潜在的な問題、または潜在的領域、あるいはこれらすべてを認識する。(i) テナント識別子のセットに対応するノードのセット、コネクタ・サービスのセット、またはディレクトリのセット、あるいはその組み合わせを管理することが、リソースにおける指数関数的増加につながる、(ii) さまざまなオペレーティング・システムが、ノードのセット、コネクタ・サービスのセット、またはディレクトリのセット、あるいはその組み合わせを多数の方法で処理する、(iii) 一部の分散ファイル・システム (`D F S : distributed file systems`) が、ポータブル・オペレーティング・システム・インターフェイス (`P O S I X : portable operating system interface`) 互換ではない、(iv) 一部の `D F S` をマウントできない、または (v) 超集中型インフラストラクチャが、リソースの使用量を減らそうとする、あるいはその組み合わせ。テナント識別子のセットに対応するノードのセット、コネクタ・サービスのセット、またはディレクトリのセット、あるいはその組み合わせを管理する従来の手段では、各テナント識別子に対応する個別のノードおよび個別のディレクトリが必要になる。

【 0 0 4 6 】

図 4 は、本発明に従う方法を表すフローチャート 4 0 0 を示している。処理が動作 `S 4`

05で開始し、マルチテナント構成サブシステムが、I/O要求をHadoopコンテナ・インスタンスから受信する。処理が動作S410に進み、マルチテナント構成サブシステムが、Hadoopコンテナ・インスタンスのテナント識別子のセットを分離する。処理が動作S415に進み、マルチテナント構成サブシステムが、少なくとも一部において、テナント識別子のセットに基づいて、Hadoopコンテナ・インスタンスを認識する。処理が動作S420に進み、マルチテナント構成サブシステムが、Hadoopコンテナ・インスタンスの許可のセットをチェックする。処理が動作S425で終了し、マルチテナント構成サブシステムが、I/O要求を処理する。

【0047】

図5は、Hadoopインスタンス502、Hadoopインスタンス504、Hadoopインスタンス506、コネクタ・サービス508、分散ファイル・システム510、および物理ノード512を含むシステム500の機能ブロック図を示している。Hadoopインスタンス502、Hadoopインスタンス504、およびHadoopインスタンス506の各々と、分散ファイル・システム510との間の通信は、コネクタ・サービス508を通して横断する。分散ファイル・システム510は、物理ノード512上に存在することによって、コネクタ・サービス508を介してすべての通信を処理できる。

【0048】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)DFSインスタンス・データのセットを分離すること、(ii)Hadoopインスタンス・データのセットを分離すること、(iii)マルチテナント認識モジュールをDFSコネクタ・サービスに導入すること、または(iv)マルチテナント能力を超集中型DFSに提供すること、あるいはその組み合わせ。超集中型DFSは、マルチテナントDFSと呼ばれることもある。本発明の一部の実施形態では、マルチテナント認識モジュールは、図4の動作S410およびS415を組み込む。他の実施形態では、図5のコネクタ・サービス508は、図4の動作S410または動作S415あるいはその両方を実行する。さらに別の実施形態では、マルチテナント構成サブシステムは、コネクタ・サービスおよび物理ノードを1対1の関係で提供する。代替の実施形態では、マルチテナント構成サブシステムは、プライベート・ネットワーク・アドレスのセットを使用してDFSインスタンスのセットを構成する。代替として、マルチテナント構成サブシステムは、プライベート・ネットワーク・アドレスを使用してDFSインスタンスのセットを構成する。一部の実施形態では、マルチテナント構成サブシステムは、ディレクトリ内のDFSインスタンスを分離する。さらに別の実施形態では、マルチテナント構成サブシステムは、少なくとも一部において、テナントに基づいてディレクトリ内のDFSインスタンスを分離する。他の実施形態では、マルチテナント構成サブシステムは、ディレクトリ内のDFSインスタンスの動作のセットを分離する。

【0049】

図6は、2つのテーブルを示している。図6の第1のテーブルは、インスタンス・コンテナ・マッピング・リスト(instance container mapping list)である。3つのコンテナを含む2つのインスタンスが示されており、そのため、6つのテナントIDが含まれている。これらの6つのテナントIDは、すべて1つのノードにマッピングされている。図6の第2のテーブルは、逆インスタンス・コンテナ・マッピング・リスト(reverse instance container mapping list)である。同じ6つのテナントIDが示されている。ただし第2のテーブルは、対応するインスタンスを決定するために分類される。

【0050】

図7は、本発明に従う方法を表すフローチャート700を示している。処理が動作S705で開始し、マルチテナント構成サブシステムが、I/O読み取り/書き込み要求をコンテナ内のHadoopジョブから受信する。処理が動作S710に進み、マルチテナント構成サブシステムが、コンテナのIPアドレスをI/O要求から取得する。処理が動作S715に進み、マルチテナント構成サブシステムが、物理ノードのIPアドレスを取得する。処理が動作S720に進み、マルチテナント構成サブシステムが、コンテナのIP

10

20

30

40

50

およびノードのIPに基づいて、インスタンス・コンテナ・マッピング・リストを照会する。処理が動作S725に進み、マルチテナント構成サブシステムが、インスタンスIPを取得する。処理が動作S730に進み、マルチテナント構成サブシステムが、インスタンスのディレクトリを取得する。処理が動作S735に進み、マルチテナント構成サブシステムが、I/O経路のセットを変換する。処理が動作S740で終了し、マルチテナント構成サブシステムが、I/O要求のセットを処理する。

【0051】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)DFSが、さまざまなホストからのファイルのセットへのアクセスを可能にするか、(ii)DFSが、ユーザのセットが、デバイスのセットにわたってファイルのセットを共有できるようにするか、または(iii)DFSが一般的なストレージ・システムであるか、あるいはその組み合わせである。DFSの例としては、IBM General Parallel File System(「GPFS(TM)」)File Placement Optimizer(「FPO」)、Red Hat Linux(R)、GlusterFS、Lustre、Ceph、およびApache Hadoop Distributed File System(「HDFS」)が挙げられる。IBMおよびGPFSは、世界中の多くの管轄区域で登録されている、International Business Machines Corporationの商標である。Linuxは、米国またはその他の国あるいはその両方における、Linus Torvaldsの登録商標である。

【0052】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)DFSをマウントすること、(ii)データをDFSから読み取ること、(iii)データをDFSに書き込むこと、(iv)POSIXアプリケーションを使用してデータをDFSから読み取ること、(v)POSIXアプリケーションを使用してデータをDFSに書き込むこと、(vi)DFSのエコシステム内でPOSIXアプリケーションを使用してデータをDFSから読み取ること、または(vii)DFSのエコシステム内でPOSIXアプリケーションを使用してデータをDFSに書き込むこと、あるいはその組み合わせ。本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)少なくとも一部において、ユーザIDに基づいて許可のセットを決定すること、(ii)少なくとも一部において、グループIDに基づいて許可のセットを決定すること、(iii)動作環境に関する許可のセットを決定すること、または(iv)オペレーティング・システムに関する許可のセットを決定すること、あるいはその組み合わせ。

【0053】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)POSIXアプリケーションを使用してDFSを実行すること、(ii)単一のコネクタ・サービスを經由してファイルのセットを転送すること、(iii)POSIXアプリケーションを使用してDFS上で単一のコネクタ・サービスを經由してファイルのセットを転送すること、または(iv)POSIXアプリケーションを使用して超集中型DFSを実行すること、あるいはその組み合わせ。本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)非POSIXアプリケーションを使用してDFSを実行すること、(ii)単一のコネクタ・サービスを經由してファイルのセットを転送すること、(iii)非POSIXアプリケーションを使用してDFS上で単一のコネクタ・サービスを經由してファイルのセットを転送すること、または(iv)非POSIXアプリケーションを使用して超集中型DFSを実行すること、あるいはその組み合わせ。本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i)DFSインスタンスのセットのクラスタのセットを作成すること、(ii)ユーザのセットのDFSインスタンスのセットのクラスタの

セットを作成すること、(i i i) ネットワーク・アドレスのセットをクラスタのセットに割り当てること、(i v) テナント識別子のセットをクラスタのセットに割り当てること、(v) ネットワーク・アドレスのセットをクラスタのセットに割り当てることであって、ネットワーク・アドレスのセットが D F S に関連していない、こと、または(v i) テナント識別子のセットをクラスタのセットに割り当てることであって、ネットワーク・アドレスのセットが D F S に関連していない、こと、あるいはその組み合わせ。

【 0 0 5 4 】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i) コネクタ・サービスの数を減らすこと、(i i) 単一のコネクタ・サービスを使用すること、(i i i) マルチテナント構成を維持するために必要なコネクタ・サービスの数を減らすこと、(i v) 指数関数的レベルでマルチテナント構成を維持するために必要なコネクタ・サービスの数を減らすこと、(v) D F S 上のクライアントの数に対応するテナント識別子の数を減らすこと、または(v i) D F S 上のクライアントの数に対応する I P アドレスの数を減らすこと、あるいはその組み合わせ。

【 0 0 5 5 】

本発明の一部の実施形態では、マルチテナント構成サブシステムがテナントの D F S クラスタを生成する。さらに別の実施形態では、マルチテナント構成サブシステムは、D F S クラスタに対応するテナント I D を生成する。D F S クラスタは、複数の要求元または複数のテナントあるいはその両方を含む第1の分散ファイル・システムと呼ばれることもある。これらの実施形態の一部では、マルチテナント構成サブシステムは、テナント I D をノードに割り当てる。

【 0 0 5 6 】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i) D F S 内のディレクトリのセットを構成すること、(i i) D F S 内のディレクトリのセットを構成し、コネクタ・サービスを再開すること、(i i i) D F S インスタンスのソフトウェア・ライブラリ・フレームワーク・インスタンスのセットを作成すること、(i v) テナント情報のセットを超集中型 D F S 内のディレクトリに格納すること、(v) 再開せずに D F S のディレクトリを認識すること、(v i) 新しい D F S インスタンスを作成せずに D F S を再開すること、(v i i) D F S クラスタをテナントに提供すること、(v i i i) テナントの D F S クラスタを維持すること、または(i x) 少なくとも一部において、ハードウェア・リソースのセットに基づいて D F S を分離すること、あるいはその組み合わせ。本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i) ソフトウェア・ライブラリ・フレームワークを構築するときに、ユーザ I D を生成すること、(i i) ソフトウェア・ライブラリ・フレームワークをコンパイルするときに、ユーザ I D を生成すること、(i i i) ソフトウェア・ライブラリ・フレームワークを構築するときに、グループ I D を生成すること、または(i v) ソフトウェア・ライブラリ・フレームワークをコンパイルするときに、グループ I D を生成すること、あるいはその組み合わせ。

【 0 0 5 7 】

本発明の一部の実施形態は、次の機能、特徴、または長所、あるいはその組み合わせのうちの1つまたは複数を含んでよい。(i) 超集中型ビッグ・データ D F S を管理すること、(i i) マルチテナント・ビッグ・データ D F S を管理すること、(i i i) クラウド・システム内の超集中型 D F S を管理すること、または(i v) 仮想システム内の超集中型 D F S を管理すること、あるいはその組み合わせ。

【 0 0 5 8 】

I V . 定義

「本発明」は、説明された対象が、出願時の初期の一連の請求項によって、審査中に書かれた修正された一連の請求項によって、または特許審査によって許可され、交付済み特

10

20

30

40

50

許に含まれる最終的な一連の請求項によって、あるいはその組み合わせによって、カバーされるということの、絶対的な指示または暗示あるいはその両方を生み出さない。「本発明」という用語は、最先端の技術を超える1つまたは複数の進歩を含んでいる可能性のある本開示の1つまたは複数の部分を示すことを補助するために使用される。「本発明」という用語およびその指示または暗示あるいはその両方についてのこの理解は、一時的かつ暫定的であり、関連する情報が開発されたとき、および請求項が修正されたときに、特許審査の過程で変更される可能性がある。

【0059】

「実施形態」については、「本発明」の定義を参照する。

【0060】

「または～あるいはその組み合わせ (And/or)」は、包含的離接であり、論理離接とも呼ばれ、一般に「包含的論理和」と呼ばれる。例えば、「A、B、またはC、あるいはその組み合わせ」という語句は、AまたはBまたはCのうちの少なくとも1つが真であるということの意味しており、「A、B、またはC、あるいはその組み合わせ」は、AおよびBおよびCの各々が偽である場合にのみ、偽になる。

【0061】

項目「のセット」は、1つまたは複数の項目が存在しており、少なくとも1つの項目が存在しなければならないが、2つ、3つ、またはそれ以上の項目が存在する可能性もあるということの意味する。項目「のサブセット」は、共通の特徴を含んでいる項目のグループ内に1つまたは複数の項目が存在するということの意味する。

【0062】

「複数の」項目は、2つ以上の項目が存在しており、少なくとも2つの項目が存在しなければならないが、3つ、4つ、またはそれ以上の項目が存在する可能性もあるということの意味する。

【0063】

「含む (Includes)」およびその変形 (例えば、含んでいる (including)、含む (include) など) は、特に明示的に示されない限り、「～を含むが、必ずしもそれに限定されない」ということを意味する。

【0064】

「ユーザ」または「サブスクライバ」は、(i) 1人の個人、(ii) 1人の個人または2人以上の人間の代わりに行動するための十分な知能を有する人工知能の実体、(iii) 1人の個人または2人以上の人間によって活動が行われている企業実体、または(iv) 1つの「ユーザ」または「サブスクライバ」として行動している任意の1つまたは複数の関連する「ユーザ」または「サブスクライバ」の組み合わせ、あるいはその組み合わせを含むが、必ずしもこれらに限定されない。

【0065】

「受信する」、「提供する」、「送信する」、「入力する」、「出力する」、および「報告する」という用語は、特に明示的に指定されない限り、(i) 物体と対象の間の関係に関する直接性の特定の程度、または(ii) 物体と対象の間に挿入された中間的なコンポーネント、中間的な動作、もしくは物、またはその組み合わせのセットの存在もしくは不在、あるいはその両方を示しているか、または暗示していると受け取られるべきではない。

【0066】

「モジュール」は、ハードウェア、ファームウェア、またはソフトウェア、あるいはその組み合わせの任意のセットであり、モジュールが次のいずれの状態であるかに関わらず、機能を実行するように動作可能である。(i) 1つにまとまって局所的に近接している、(ii) 広い領域にわたって分散している、(iii) 大きい一群のソフトウェア・コード内で1つにまとまって近接している、(iv) 一群のソフトウェア・コード内に存在する、(v) 1つのストレージ・デバイス、メモリ、または媒体内に存在する、(vi) 機械的に接続されている、(vii) 電氣的に接続されている、または(viii) デー

10

20

30

40

50

タ通信で接続されている、あるいはその組み合わせの状態にある。「サブモジュール」は、「モジュール」内の「モジュール」である。

【0067】

「コンピュータ」は、大きいデータ処理能力または機械可読命令読み取り能力あるいはその両方を備える任意のデバイスであり、デスクトップ・コンピュータ、メインフレーム・コンピュータ、ラップトップ・コンピュータ、フィールドプログラマブル・ゲート・アレイ（FPGA）ベース・デバイス、スマートフォン、パーソナル・デジタル・アシスタント（PDA）、人体に装着または挿入されるコンピュータ、組み込みデバイス型コンピュータ、または特定用途向け集積回路（ASIC：application-specific integrated circuit）ベース・デバイス、あるいはその組み合わせを含むが、必ずしもこれらに限定されない。

10

【0068】

「電氣的に接続される」とは、介在する要素が存在するように間接的に電氣的に接続されるか、または直接的に電氣的に接続されることを意味する。「電気接続」は、コンデンサ、インダクタ、変圧器、真空管などの要素を含んでよいが、これらに限定される必要はない。

【0069】

「機械的に接続される」とは、中間的なコンポーネントを介して行われる間接的な機械的接続または直接的な機械的接続のいずれかを意味する。「機械的に接続される」は、固定された機械的接続、および機械的に接続されコンポーネント間の相対運動を許容する機械的接続を含む。「機械的に接続される」は、溶接接続、半田接続、締め具（例えば、くぎ、ボルト、ねじ、ナット、面ファスナ、結び目、リベット、簡易脱着接続、掛け金、または磁氣的接続、あるいはその組み合わせ）による接続、圧力ばめ接続、摩擦適合接続、重力によって引き起こされるかみ合いによって固定される接続、旋回可能もしくは回転可能な接続、またはスライド可能な機械的接続、あるいはその組み合わせを含むが、これらに限定されない。

20

【0070】

「データ通信」は、現在知られているか、または将来開発される任意の種類のデータ通信方式を含むが、必ずしもこれらに限定されない。「データ通信」は、無線通信、有線通信、または無線部分および有線部分を含む通信経路、あるいはその組み合わせを含むが、必ずしもこれらに限定されない。「データ通信」は、（i）直接データ通信、（ii）間接データ通信、または（iii）形式、パケット化の状態、媒体、暗号化の状態、もしくはプロトコル、またはその組み合わせが、データ通信の過程全体にわたって一定のままであるデータ通信、あるいはその組み合わせに、必ずしも限定されない。

30

【0071】

「実質的に人間が介入しない」という語句は、人間による入力がほとんど、または全く存在せずに（多くの場合、ソフトウェアなどの機械論理の動作によって）自動的に発生するプロセスを意味する。「実質的に人間が介入しない」ことを含む例としては、（i）コンピュータが複雑な処理を実行しており、電力系統の停止に起因して処理の継続が中断されないように、人間がコンピュータを代替電源に切り替えること、（ii）コンピュータが、リソースを大量に使用する処理を実行しようとしており、人間が、リソースを大量に使用する処理が本当に実行されるべきかどうかを確認すること（この場合、分離していると考えられる確認のプロセスは、実質的な人間の介入を伴うが、人間によって行われる必要のある単純なはい/いいえ形式の確認にもかかわらず、リソースを大量に使用する処理は、どのような実質的な人間の介入も含まない）、および（iii）コンピュータが、機械論理を使用して重要な決定（例えば、悪天候を予想して、すべての飛行機を飛行禁止にすることの決定）を行うが、重要な決定を実施する前に、コンピュータが、単純なはい/いいえ形式の確認を人間側から得なければならないこと、が挙げられる。

40

【0072】

「自動的に」とは、「人間による介入なしで」ということを意味する。

50

【 0 0 7 3 】

「リアルタイム」（および形容詞の「リアルタイムの」）という用語は、前述した情報処理で妥当な応答時間を実現するために十分短い期間である、任意の時間枠を含む。さらに、「リアルタイム」（および形容詞の「リアルタイムの」）という用語は、一般に「ほぼリアルタイム」と呼ばれる、前述したオンデマンドの情報処理で妥当な応答時間を実現するために十分短い期間である、任意の（例えば、1秒未満または数秒以内の）時間枠を通常は含む。これらの用語は、正確に定義するのは困難であるが、当業者によって十分理解されている。

【 符号の説明 】

【 0 0 7 4 】

- 1 0 0 ネットワーク・コンピュータ・システム
- 1 0 2 マルチテナント構成サブシステム
- 1 0 4 ユーザ・サブシステム
- 1 0 6 , 1 0 8 仮想コンテナ・サブシステム
- 1 1 0 ソフトウェア・ライブラリ・フレームワーク
- 1 1 1 コネクタ・サービス
- 1 1 4 通信ネットワーク
- 2 0 0 マルチテナント構成コンピュータ
- 2 0 1 通信ユニット
- 2 0 4 プロセッサ・ユニット
- 2 0 6 I / O インターフェイス・セット
- 2 0 8 メモリ・デバイス
- 2 1 0 永続的ストレージ・デバイス
- 2 1 2 ディスプレイ・デバイス
- 2 1 4 外部デバイス
- 2 1 6 R A M デバイス
- 2 1 8 キャッシュ・メモリ・デバイス
- 3 0 0 マルチテナント構成プログラム

10

20

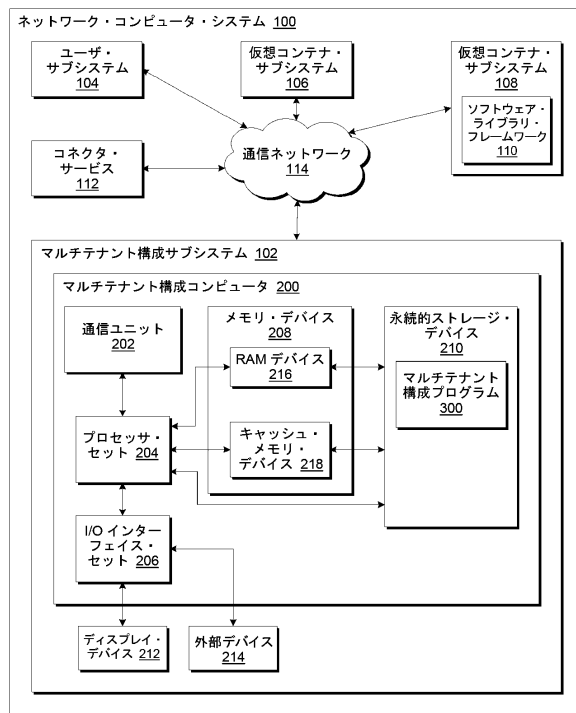
30

40

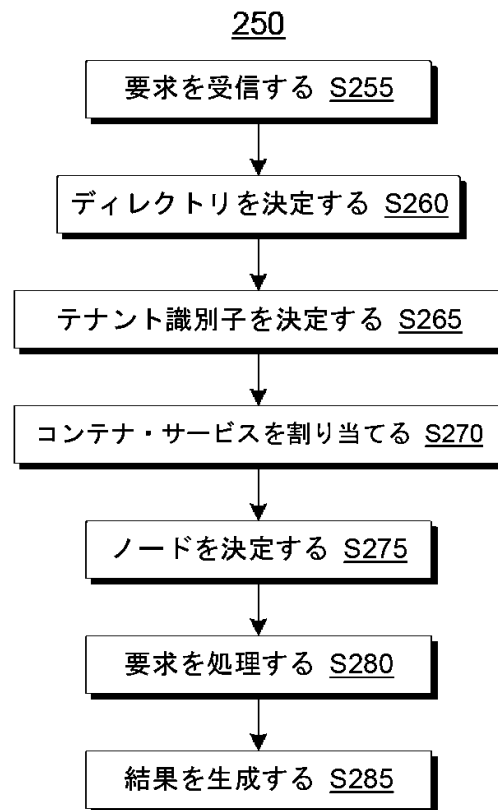
50

【図面】

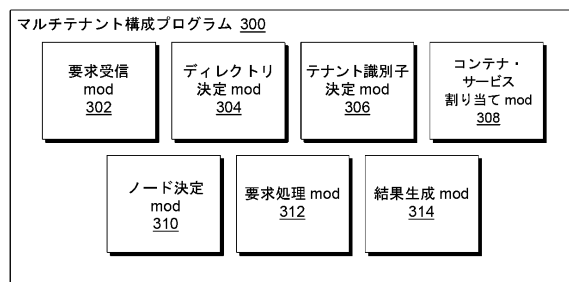
【図 1】



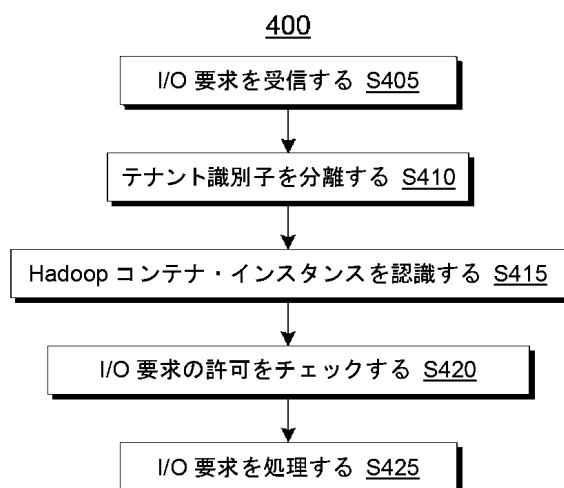
【図 2】



【図 3】



【図 4】



10

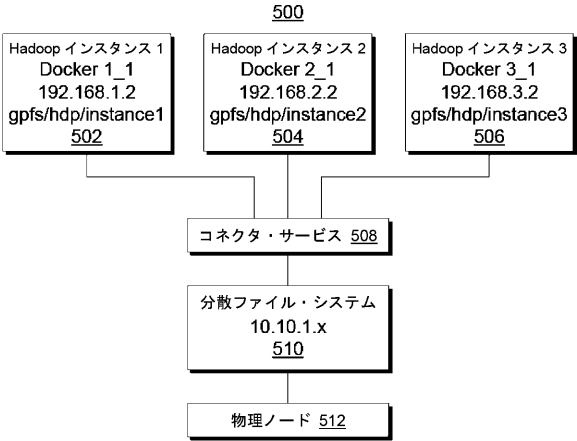
20

30

40

50

【図 5】

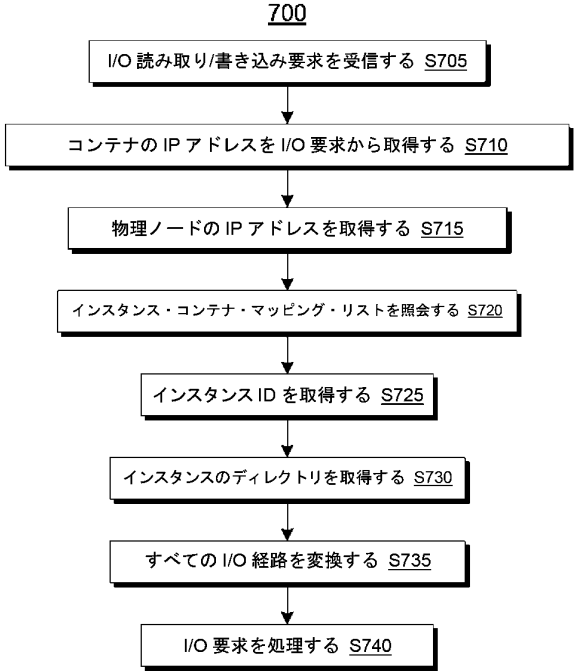


【図 6】

インスタンス	コンテナ	テナント ID	ノード
インスタンス 1	コンテナ 1_1	192.168.1.2	ノード 1
インスタンス 1	コンテナ 2_1	192.168.1.3	ノード 1
インスタンス 1	コンテナ 3_1	192.168.1.4	ノード 1
インスタンス 2	コンテナ 1_2	192.168.2.2	ノード 1
インスタンス 2	コンテナ 2_2	192.168.2.3	ノード 1
インスタンス 2	コンテナ 3_2	192.168.2.4	ノード 1
ノード	コンテナ	テナント ID	インスタンス
ノード 1	コンテナ 1_1	192.168.1.2	インスタンス 1
ノード 1	コンテナ 2_1	192.168.1.3	インスタンス 1
ノード 1	コンテナ 3_1	192.168.1.4	インスタンス 1
ノード 1	コンテナ 1_2	192.168.2.2	インスタンス 2
ノード 1	コンテナ 2_2	192.168.2.3	インスタンス 2
ノード 1	コンテナ 3_2	192.168.2.4	インスタンス 2

10

【図 7】



20

30

40

50

フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

(72)発明者 チェン、ヨン

中華人民共和国 1 1 1 0 0 1 9 3 ベイジン ドンベイワン ウエスト・ロード 8 ジョングァンツ
ン・ソフトウェア・パーク 2 8

(72)発明者 ユアン、チェン、カイ

中華人民共和国 1 1 1 0 0 1 9 3 ベイジン ドンベイワン ウエスト・ロード 8 ジョングァンツ
ン・ソフトウェア・パーク 2 8

(72)発明者 フォン、ティエン

中華人民共和国 1 0 0 0 1 9 3 ベイジン チャン・ピン・ディストリクト 1 5 - 1 - 3 0 2 ロ
ン・ユエ・ヤン 4

(72)発明者 ワン、シン

中華人民共和国 1 1 1 0 0 1 9 3 ベイジン ドンベイワン ウエスト・ロード 8 ジョングァンツ
ン・ソフトウェア・パーク 2 8

(72)発明者 バオ、シャオ、ミン

中華人民共和国 1 1 1 0 0 1 9 3 ベイジン ドンベイワン ウエスト・ロード 8 ジョングァンツ
ン・ソフトウェア・パーク 2 8

審査官 北村 学

(56)参考文献 特開 2 0 0 0 - 1 4 8 5 6 5 (J P , A)

特開 2 0 0 7 - 1 8 8 2 0 9 (J P , A)

特開 2 0 0 7 - 2 8 7 1 8 0 (J P , A)

特開 2 0 1 1 - 1 3 4 0 3 7 (J P , A)

特表 2 0 1 4 - 5 1 7 9 4 1 (J P , A)

特開 2 0 1 5 - 2 1 9 8 5 2 (J P , A)

米国特許第 0 9 0 6 9 7 7 8 (U S , B 1)

米国特許出願公開第 2 0 1 4 / 0 0 0 6 7 0 8 (U S , A 1)

杵淵 雄樹, 検証解剖! ジャーナリングファイルシステム E x t 3 , U N I X U S E R , 日
本, ソフトバンクパブリッシング株式会社, 2005年07月01日, 第14巻 第7号, pp. 62-68
水吉 俊幸, L i n u x ビジネス・ソリューション 第3回, 日経バイト, 日本, 日経 B P 社
, 1999年11月22日, 第198号, pp. 182 ~ 189Apache Hadoop 2.7.2 - Introduction , The Apache Software Foundation , 2016年03月11日
, [https://web.archive.org/web/20160311231004/https://hadoop.apache.org/docs/curre](https://web.archive.org/web/20160311231004/https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/filesystem/introduction.html)
[nt/hadoop-project-dist/hadoop-common/filesystem/introduction.html](https://web.archive.org/web/20160311231004/https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/filesystem/introduction.html)

(58)調査した分野 (Int.Cl. , D B 名)

I P C G 0 6 F 1 6 / 0 0 - 1 6 / 9 5 8

G 0 6 F 9 / 5 0