

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4263421号
(P4263421)

(45) 発行日 平成21年5月13日(2009.5.13)

(24) 登録日 平成21年2月20日(2009.2.20)

(51) Int.Cl. F I
G06F 12/00 (2006.01)
 G06F 12/00 545A
 G06F 12/00 531D
 G06F 12/00 535B

請求項の数 10 (全 28 頁)

<p>(21) 出願番号 特願2002-86870 (P2002-86870) (22) 出願日 平成14年3月26日(2002.3.26) (65) 公開番号 特開2002-358226 (P2002-358226A) (43) 公開日 平成14年12月13日(2002.12.13) 審査請求日 平成17年3月23日(2005.3.23) (31) 優先権主張番号 60/278,905 (32) 優先日 平成13年3月26日(2001.3.26) (33) 優先権主張国 米国 (US) (31) 優先権主張番号 10/005,629 (32) 優先日 平成13年12月5日(2001.12.5) (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 500046438 マイクロソフト コーポレーション アメリカ合衆国 ワシントン州 9805 2-6399 レッドモンド ワン マイ クロソフト ウェイ (74) 代理人 100077481 弁理士 谷 義一 (74) 代理人 100088915 弁理士 阿部 和夫 (72) 発明者 アトゥル アジャ アメリカ合衆国 98007 ワシントン 州 ベルビュー ノースイースト 35 ストリート 14590 ナンバーディー 204</p>
---	--

最終頁に続く

(54) 【発明の名称】 サーバレス分散ファイルシステム

(57) 【特許請求の範囲】

【請求項1】

ファイルを記憶するのに使用する階層的な名前空間のサブツリーを管理する1または複数のコンピュータからなるディレクトリグループを決定する方法であって、

第1のディレクトリを管理する第1のディレクトリグループのコンピュータが、前記第1のディレクトリのサブツリーの一部の管理を委任する第2のディレクトリグループのコンピュータを選択するステップと、

前記第1のディレクトリグループのコンピュータが、前記サブツリーの委任証明書を生成するステップであって、前記委任証明書は、

第3のディレクトリグループのコンピュータが、前記サブツリーの名前空間ルートの管理に責任を負うことを証明する、第1のデジタル署名された証明書と、

前記第2のディレクトリグループのコンピュータの権限を、前記名前空間ルートまでさかのぼって、前記第3のディレクトリグループのコンピュータの証明書と関連付けることを可能にする、第2のデジタル署名された証明書とを含む、ステップと、

前記第1のディレクトリグループのコンピュータが、前記委任証明書にデジタル署名するステップと、

前記第1のディレクトリグループのコンピュータが、前記第2のディレクトリグループのコンピュータに前記委任証明書を発行するステップと

を備えたことを特徴とする方法。

【請求項2】

10

20

前記委任証明書にデジタル署名することは、前記委任証明書に、複数のコンピュータがデジタル署名させることを備えたことを特徴とする請求項 1 に記載の方法。

【請求項 3】

前記第 2 の ディレクトリグループ のコンピュータは、ビザンチンフォルトトレラントグループを構成すること を特徴とする請求項 1 に記載の方法。

【請求項 4】

前記第 1 および第 2 のデジタル署名された証明書は、
委任されるサブツリーのルートの識別と、
前記委任されるサブツリーの識別と、
前記サブツリーを管理するディレクトリグループのメンバの識別と
を含むこと を特徴とする請求項 1 に記載の方法。

10

【請求項 5】

前記第 2 の ディレクトリグループ のコンピュータは、前記第 1 の ディレクトリグループ のコンピュータに 含まれていること を特徴とする請求項 4 に記載の方法。

【請求項 6】

前記第 1 のデジタル署名された証明書は、認証局 (CA) によってデジタル署名された証明書であることを特徴とする請求項 1 に記載の方法。

【請求項 7】

前記委任証明書は、前記第 2 のデジタル署名された証明書から前記第 1 のデジタル署名された証明書への、証明書のチェーンの確立を可能にする、デジタル署名された 1
 または複数の追加の証明書を含むことを特徴とする請求項 1 に記載の方法。

20

【請求項 8】

サーバレス分散ファイルシステムにおいて、ファイルを記憶するのに使用する階層的な名前空間のサブツリーを管理する 1 または複数のコンピュータからなる ディレクトリグループ を決定する、第 1 のディレクトリを管理する第 1 のディレクトリグループのコンピュータ であって、

プロセッサと、

該プロセッサに結合されたメモリとを備え、

該メモリに格納されたプログラム命令を、前記プロセッサが実行すると、

前記第 1 の ディレクトリのサブツリーの一部の管理を委任する第 2 のディレクトリグループ のコンピュータを識別するステップと、

30

前記サブツリーの委任証明書を生成するステップであって、前記委任証明書は、
第 3 のディレクトリグループのコンピュータが、前記サブツリーの名前空間ルートの管理に責任を負うこと を証明する、第 1 のデジタル署名された証明書と、

前記第 2 の ディレクトリグループ のコンピュータの権限を、前記名前空間ルートまでさかのぼって、前記第 3 の ディレクトリグループ のコンピュータの証明書と関連付けることを可能にする、第 2 のデジタル署名された証明書とを含む、ステップと、

前記委任証明書にデジタル署名するステップと、

前記第 2 の ディレクトリグループ のコンピュータに前記委任証明書を発行するステップと

40

を実行することを特徴とするコンピュータ。

【請求項 9】

前記第 2 の ディレクトリグループ のコンピュータは、ビザンチンフォルトトレラントグループを構成すること を特徴とする請求項 8 に記載のコンピュータ。

【請求項 10】

前記第 1 および第 2 のデジタル署名された証明書は、

委任されるサブツリーのルートの識別と、

前記委任されるサブツリーの識別と、

前記サブツリーを管理するディレクトリグループのメンバの識別と

を含むこと を特徴とする請求項 8 に記載のコンピュータ。

50

【発明の詳細な説明】**【0001】****【発明の属する技術分野】**

本発明は、コンピュータネットワークおよびファイルシステムに関し、より詳細には、サーバレス分散ファイルシステム (a serverless distributed file system) に関する。

【0002】**【従来の技術】**

ファイルシステムは、コンピュータシステム上に記憶されたファイルおよび他のデータオブジェクトを管理する。ファイルシステムは元来、コンピュータのオペレーティングシステム中に構築され、常駐する記憶媒体上にローカルに記憶されたファイルへのアクセスを容易にする。コンピュータがネットワーク化されると、いくつかのファイル記憶機能は、個々のユーザマシンから専用の記憶サーバへオフロードされ、これらのユーザマシンに代って専用記憶サーバが多数のファイルを記憶する。ファイルが必要になると、ユーザマシンは、単にサーバにファイルを要求した。このサーバベースのアーキテクチャにおいて、ネットワークを介して遠隔に、記憶サーバに記憶されたファイルの管理およびアクセスを容易にするようにファイルシステムが拡張された。

10

【0003】**【発明が解決しようとする課題】**

今日、ファイル記憶は、中央記憶サーバ上にファイルを記憶するよりも、ネットワーク化されたさまざまなコンピュータ上にファイルを記憶するモデルへと移行している。サーバのないアーキテクチャは、ファイルシステムに新たな課題を提起している。具体的な1つの課題は、多くの異なるコンピュータ上に分散したファイルを、コンピュータのいくつかを任意の時刻にアクセス不能にするにもかかわらず、高い信頼性で記憶し、アクセスできる方法で管理し、同時に、無許可ユーザによるファイルへのアクセスを防止することに関する。

20

【0004】

本発明は、サーバレス分散ファイルシステムのこれらの課題に取り組み、有効な解決方法を提供する。

【0005】**【課題を解決するための手段】**

本明細書には、サーバレス分散ファイルシステムを開示する。

30

【0006】

一態様によれば、ファイルとディレクトリは、サーバレス分散ファイルシステム内で異なる方法で管理される。ディレクトリは、ビザンチンフォルトトレラントグループ (Byzantine-fault-tolerant groups) を使用して管理され、一方、ファイルは、ビザンチンフォルトトレラントグループを使用せずに管理される。これは、対応するディレクトリエントリのコピーよりも少ない数のファイルコピーを記憶する結果、性能が向上する。

【0007】

他の態様によれば、ファイルシステムは、階層的名前空間を使用してファイルを記憶する。ファイルは、複数のコンピュータを横断して分散され、これら複数のコンピュータの各々は、クライアントコンピュータおよびサーバコンピュータとして動作することができ、他のコンピュータを信用する必要がない。

40

【0008】

他の態様によれば、ファイルシステム内の1または複数のディレクトリを管理する責任が、ディレクトリグループに割り当てられる。ディレクトリグループの各メンバは、システムに参加したコンピュータであり、ディレクトリグループは、各ディレクトリ中のオブジェクト (例えばファイルおよびディレクトリ) へのアクセスを制御する複数のロック (locks) を使用することができる。ロックは、オブジェクトのオープンを制御する第1のロックセット、およびオブジェクト中のデータへのアクセスを制御する第2のロックセットを含む。

50

【 0 0 0 9 】

【 発明の実施の形態 】

本明細書を通して、同じ構成要素および/または特徴を参照するのに同じ符号を使用した。

【 0 0 1 0 】

以下の議論は、ネットワーク化された複数のコンピュータ上で稼動し、中央サーバまたはサーバ群にファイルを記憶するのではなく、これらのネットワーク化されたコンピュータ上にファイルを記憶する、共生 (symbiotic)、サーバレス (serverless)、分散ファイルシステムを対象とする。共生とは、マシンが相互に協力はするが、完全には信用しないことを意味する。このファイルシステムは、ディスク記憶装置を直接には管理せず、オペレーティングシステムに組み込まれたファイルシステム (例えば、Windows (登録商標) NTのファイルシステム) などの、ローカルマシン上に存在するファイルシステムに依存する。

10

【 0 0 1 1 】

本明細書の議論では、公開鍵、暗号化およびデジタル署名に言及する。公開鍵暗号とは一般に、ひとまとめに鍵対と呼ばれる公開鍵と秘密鍵の使用を指す。実体 (例えば、ユーザ、ソフトウェアアプリケーションなど) は、秘密鍵を秘密にしておき、公開鍵を他の実体に公開する。データは一般に平文と呼ばれ、暗号化アルゴリズムおよび公開鍵を使用して暗号化することができる。この場合、暗号化した結果 (一般に暗号文と呼ばれる) は、対応する秘密鍵を知らなければ簡単には復号できないが、対応する秘密鍵を知っていれば比較的簡単に復号することができる。同様に、暗号化アルゴリズムおよび秘密鍵を使用してデータにデジタル署名することができる。この場合、署名は、対応する公開鍵を使用して簡単に検証できるが、秘密鍵がなければ簡単に生成することができない。本明細書の議論は、読者が暗号法の基本的な理解を有することを前提とする。暗号法の基本については、Bruce Schneier著「Applied Cryptography :Protocols, algorithms, and Source Code in C」、John Wiley & Sons社刊、1994年 (または第2版、1996年) を参照されたい。

20

【 0 0 1 2 】

(サーバレス分散ファイルシステム)

図1に、サーバレス分散ファイルシステムをサポートする例示的なネットワーク環境100を示す。4台のクライアントコンピューティング装置102、104、106および108が、データ通信ネットワーク110を介して相互に結合されている。4台のコンピューティング装置を図示したが、ネットワーク環境100には、さまざまな数 (4台を越えるまたは4台未満) の装置を含めることができる。

30

【 0 0 1 3 】

ネットワーク110は、多種多様なデータ通信ネットワークのうちの1つを表す。ネットワーク110は、公的部分 (例えばインターネット)、私的部分 (例えば企業内ローカルエリアネットワーク (LAN))、あるいは公的部分と私的部分との組合せを含むことができる。ネットワーク110は、有線および無線媒体を含む従来のさまざまな通信媒体のうち、1または複数の媒体を使用して実現することができる。公衆のおよび私有のプロトコルを含むさまざまな通信プロトコルのいずれかを使用して、ネットワーク110経由でデータを伝達することができる。このようなプロトコルの例にはTCP/IP、IPX/SPX、NetBEUIなどがある。

40

【 0 0 1 4 】

コンピューティング装置102~108は、広範囲にわたるコンピューティング装置の1つを表し、互いに同じ装置であってもよいし、または異なる装置であってもよい。装置102~108は例えば、デスクトップコンピュータ、ラップトップコンピュータ、ハンドヘルドまたはポケットコンピュータ、パーソナルデジタルアシスタント (PDA)、携帯電話、インターネット機器、家庭用電子装置、ゲーム機などである。

【 0 0 1 5 】

50

2台以上の装置102～108が動作して、サーバレス分散ファイルシステムを実現する。サーバレス分散ファイルシステムに参加する実際の装置は、その時々で変更することができる。すなわち、新たな装置をシステムに追加し、他の装置をシステムから除くことができる。分散ファイルシステムを実現する（分散ファイルシステムに参加する）各装置102～106は、ローカル記憶または分散記憶として使用するよう割り当てられた大容量記憶装置（例えば、ハードディスクドライブ）の部分を持つ。ローカル記憶は、ユーザが、分散ファイルシステム構造にではなく、自分のローカルマシン上に記憶しておきたいデータに対して使用される。分散記憶部分は、その装置（または他の装置）のユーザが分散ファイルシステム構造に記憶しておきたいデータに対して使用される。

【0016】

図1に示した例では、ネットワーク110に接続されたある装置が、分散部分とローカル部分の両方を含む1または複数の大容量記憶装置を持つ。分散またはローカル記憶に割り当てられる量は、装置ごとに異なる。例えば装置102では、ローカル部分122に比べて分散システム部分120に大きな割合が割り当てられており、装置104は、ローカル部分126とほぼ同じサイズの分散システム部分124を含み、装置106では、ローカル部分130に比べて分散システム部分128に小さな割合が割り当てられている。記憶域の複数部分の分離は、記憶装置単位で実施し（例えば、1台のハードドライブを分散システムで使用するよう指定し、別のハードドライブをローカル使用専用に指定する）、および/または単一の記憶装置内で実施する（例えば、1台のハードドライブの一部分を分散システムで使用するよう指定し、別の部分をローカル使用向けに指定する）ことができる。分散記憶またはローカル記憶に割り当てられる量は、その時々で変更することができる。コンピューティング装置108などのネットワーク110に接続された他の装置が、分散ファイルシステムのいずれの部分をも実現せず、したがってその大容量記憶装置内に、分散システムが使用するよう割り振られた部分を持たないことがある。したがって装置108は、ローカル部分132だけしか持たない。

【0017】

分散ファイルシステム150は、異なるコンピューティング装置102～106に1または複数のファイルコピーを記憶するよう動作する。コンピュータのユーザが新たなファイルを作成したとき、ユーザは、自分のコンピューティング装置のローカル部分にファイルを記憶するのか、または分散ファイルシステムにファイルを記憶するのかを選択することができる。分散ファイルシステム150に記憶する場合には、ファイルが、1台または数台の装置102～106の大容量記憶装置の分散システム部分に記憶される。ファイルを作成したユーザは一般に、どの装置102～106上にファイルを記憶するかを制御することができず、そのファイルがどの装置102～106上に記憶されたかも知ることができない。さらに、複製された複数のファイルコピーが保存され、たとえファイルが保存されているコンピューティング装置102～106のうちの1台が使用できない（例えば、パワーダウンまたは誤動作しているなど）場合であっても、ユーザが続いてファイルを検索することができる。

【0018】

分散ファイルシステム150は、各装置102～106上の1または複数のコンポーネントによって実現され、これによってファイルシステムを調整する中央サーバの必要性が排除される。これらコンポーネントは、特定のファイルをどこに記憶するか、いくつのファイルコピーを作成して異なる装置上に記憶するかなどを決定するよう動作する。どの装置にどのファイルが記憶されるかは多くの因子によって決まり、これには、分散ファイルシステム内の装置数、各装置からファイルシステムに割り振られた記憶空間、保存するファイルコピーの数、暗号法上安全な乱数、装置上にすでに記憶されているファイルの数などが含まれる。したがって、この分散ファイルシステムによってユーザは、正確にどのコンピューティング装置上にファイルが記憶されているのかを一切知らなくとも、ファイル（フォルダまたはディレクトリも同様に）を作成し、アクセスすることができる。

【0019】

分散ファイルシステム 150 は、システム 150 内の多数のコンピュータをサポートするため、拡張が容易であるように設計される。システム 150 内の装置上のコンポーネントが使用するプロトコルおよびデータ構造は、システム内のコンピュータの数と釣り合うようには設計されず、これによって、多数のコンピュータに合わせて容易に拡張することができる。

【0020】

ファイルシステムによって記憶されるファイルは、各装置 102 ~ 106 に分散され、暗号化された形態で記憶される。新しいファイルが作成されると、他の装置に伝達して記憶する前に、ファイルを作成した装置がファイルを暗号化する。新しいファイルのディレクトリエントリ（ファイル名を含んでいる）も、他の装置に伝達され記憶される。他の装置は、暗号化されたファイルを記憶する装置と同じである必要はない（一般には異なる装置に記憶される）。さらに、新しいフォルダまたはディレクトリが作成された場合も、ディレクトリエントリ（フォルダ名またはディレクトリ名を含んでいる）が他の装置に伝達され記憶される。本明細書においてディレクトリエントリは、ファイルシステムのディレクトリに追加することができ、ファイル名とディレクトリ名（またはフォルダ名）の両方を含む、任意のエントリを指す。

10

【0021】

分散ファイルシステム 150 は、1台の装置 102 ~ 106 上に記憶されたデータを無許可ユーザが読めないように設計される。したがって、装置 102 が作成し、装置 104 に記憶されたファイルを、装置 104 のユーザは（許可を得ない限り）読むことができない。このようなセキュリティを実現するため、ファイルの内容、ディレクトリエントリ中の全てのファイル名およびディレクトリ名が暗号化され、許可されたユーザだけに復号鍵が与えられる。したがって、装置 102 が作成したファイルを装置 104 は記憶することができるが、装置 104 のユーザがそのファイルの許可されたユーザでない限り、ファイルの内容も、そのディレクトリエントリ中のファイル名も復号することができない（したがって読むことができない）。

20

【0022】

分散ファイルシステム 150 は、1または複数の名前空間ルートを有し、各名前空間ルートのもとに複数のサブツリーを有する階層的記憶構造を使用する。異なるサブツリーの管理を異なるコンピュータ群に委任し、これによって名前空間ルートまたは特定のサブツリーを管理するコンピュータの負担が過大になるのを防ぐことができる。

30

【0023】

分散ファイルシステム 150 はさらに、ファイルの記憶とファイルに対応するディレクトリエントリの記憶を異なる方法で管理する。システム 150 内に記憶されるファイルは複製され、システム内の異なる複数のコンピュータ上に保存される。さらに、そのファイルに対してディレクトリエントリが生成され、システム内の異なる複数のコンピュータ上に保存される。保存されるファイルコピーの数よりも多くのディレクトリエントリが保存される。一実施態様では、ディレクトリエントリが、ピザンチンフォルトトレラントグループの一部であるコンピュータ上に記憶される。これについては後に詳細に論じる。

40

【0024】

分散ファイルシステム 150 はさらに、ディレクトリおよびファイルの読取りまたは書込みを制御することができるディレクトリ/ファイルロック機構を使用する。ピザンチンググループ内のコンピュータとともに使用すると、使用されたロック機構は、ローカルに実行することができるオペレーションの数を増やすことによって、ディレクトリグループによるアクションを要求することなく、性能を向上させようとする。これについては後に詳細に論じる。

【0025】

分散ファイルシステム 150 内の全てのコンピュータ 102 ~ 106 は、3つの機能を有することができる。すなわち、ローカルユーザのクライアントとなり、システム内に記憶された暗号化されたファイルコピーの保管場所となり、1または複数のディレクトリを維

50

持する一群のコンピュータのメンバとなることができる。

【 0 0 2 6 】

一般に、コンピュータ 1 0 2 ~ 1 0 6 のユーザが所与のディレクトリ内のファイルを開くと、そのコンピュータは、ビザンチンフォルトトレラントプロトコルを使用して、そのディレクトリを集散的に管理している一組のコンピュータ（「ビザンチングループ」または「ディレクトリグループ」と呼ばれる）にリクエストを送る。ビザンチングループは、そのコンピュータにファイルロックを許可し、そのファイルにローカル更新を実施し（書込みロックの場合）、続いてそれらの更新をビザンチングループへプッシュすることを可能にする。コンピュータが最近このファイルにアクセスしていた場合、このコンピュータはおそらく、暗号化されたファイル内容のコピーをローカルキャッシュ内に有し、そのためコンピュータは、キャッシュされているコピーを検索し、それを復号するだけでよく、その後、コンピュータはファイルの読取りまたは書込みを開始することができる。コンピュータが最近に、ファイルの現行バージョンにアクセスしていない場合、コンピュータは、そのファイルを記憶したコンピュータの 1 つから暗号化されたファイルのコピーを検索する。どのコンピュータが現行コピーを保持しているかについての情報は、ロック許可とともにビザンチングループによって提供される。ファイルを記憶した 1 台または数台のコンピュータがダウンしている場合には、別のコンピュータからファイルを検索する。ビザンチングループはさらに、取り込んだファイルを確認するためコンピュータが使用する、ファイル内容の暗号ハッシュを提供する。

10

【 0 0 2 7 】

（ファイルの暗号化）

ファイルは、「収束化暗号（convergent encryption）」として知られる技術を使用して暗号化される。収束化暗号は、以下の 2 つの特性を有する。第 1 に、暗号化可能な 2 またはそれ以上のオブジェクトが同一である場合、それらを暗号化して個々の暗号オブジェクトを提供するのに異なる暗号化鍵を使用する場合でも、暗号化鍵のいくつかにアクセスし、暗号化可能なオブジェクトが同一であることを、暗号オブジェクトを調べて決定する必要がない。第 2 に、2 またはそれ以上の暗号化可能なオブジェクトが同一であり、それらを異なる暗号化鍵で暗号化する場合、全ての暗号オブジェクトを記憶するのに必要な総記憶空間は、暗号化可能な単一のオブジェクトを記憶するのに必要な空間に各々の異なる暗号化鍵に対する一定量の記憶域を加えたものに比例する。

20

30

【 0 0 2 8 】

一般に、ノンバージョン暗号化によれば、ファイル F （または他の暗号化可能オブジェクトの任意のタイプ）は、一方向ハッシュ関数 h （例えば SHA 、 $MD5$ など）を使用してハッシュされ、ハッシュ値 $h(F)$ を生成する。次いで、対称暗号（例えば $RC4$ 、 $RC2$ など）を使用し、ハッシュ値を鍵として、すなわち $E_{h(F)}(F)$ （ F ）としてファイル F を暗号化する。次に、この暗号化されたファイルへの読取りアクセスが許可された各許可ユーザに対して、読取りアクセス制御エントリを作成する。書込みアクセス制御は、そのファイルのディレクトリエントリを記憶したディレクトリサーバによって統率される。読取りアクセス制御エントリは、任意の数の鍵 K_1 、 K_2 、 \dots 、 K_m を用いて、ファイルのハッシュ値 $h(F)$ を暗号化し、 $E_{K_1}(h(F))$ 、 $E_{K_2}(h(F))$ 、 \dots 、 $E_{K_m}(h(F))$ を得ることによって、形成される。一実施態様では、各々の鍵 K が、非対称暗号（例えば RSA ）の公開 / 秘密鍵対のユーザの公開鍵である。

40

【 0 0 2 9 】

収束化暗号では、サーバレス分散ファイルシステム 1 5 0 の間で、ファイルの 1 つの暗号化バージョンが記憶され、複製される。ファイルの暗号化バージョンとともに、アクセスを有する許可ユーザ数に応じた 1 または複数のアクセス制御エントリが記憶される。したがって、分散ファイルシステム 1 5 0 内のファイルは以下の構造を有する。

[$E_{h(F)}(F)$, $\langle E_{K_1}(h(F)) \rangle$, $\langle E_{K_2}(h(F)) \rangle$, \dots , $\langle E_{K_m}(h(F)) \rangle$]

収束化暗号の 1 つの利点は、暗号化されたファイルをファイルシステムが評価して、復号に訴えることなく（したがって暗号化鍵を知らなくても）、それが他のファイルと同一か

50

どうかを決定することができる点である。許可ユーザのアクセス制御エントリを残りのファイルに追加することによって、不要な複製ファイルを除去することができる。他の利点は、おそらくギガバイトに達する暗号化されたファイルに比べて、アクセス制御エントリのサイズが非常に小さいことである。その結果、各々のファイルに記憶されるオーバーヘッド情報の量が小さくなる。これは、ファイルを記憶するために使用される総記憶空間が、暗号化された単一のファイルを記憶するのに必要な空間に、そのファイルの各追加の許可された読者に対する一定量の記憶域を加えたものに比例するという特性を可能にする。

【0030】

収束化暗号の詳細については、Douceur他の名義で2000年5月5日に出願され、Microsoft Corporationに譲渡された「Encryption Systems and Methods for Identifying and Coalescing Identical Objects Encrypted with different Keys」という名称の同時係属米国特許出願第09/565821号を参照されたい。この出願は、参照によって本明細書に組み込まれる。

10

【0031】

(ディレクトリエントリの暗号化)

ディレクトリエントリ中のファイル名およびディレクトリ名は、「排他的暗号(exclusive encryption)」と呼ばれるプロセスを使用して暗号化する。排他的暗号は、ディレクトリエントリ中のファイル名およびディレクトリ名を暗号化された形で記憶することを可能にし、これによって無許可ユーザが、ファイル名またはディレクトリ名に基づいて情報を不適切に獲得することを防ぐ。さらに、排他的暗号は、以下の3つの特性を有する。第1に、1つのディレクトリ内の暗号化されたエントリは、同じ名前に復号されることがない。第2に、1つのディレクトリ内の暗号化されたエントリが全て、構文上合法な名前に復号される。第3に、ディレクトリを維持するディレクトリグループが、エントリの平文名にアクセスすることができない。したがってファイルシステム150は、ディレクトリ内の2つのエントリが同じ名前の暗号化でないこと、およびディレクトリ内の全てのエントリが構文上合法な名前の暗号化であることを保証することができ、同時に、ディレクトリを維持する装置が、エントリの平文名にアクセスできないことを保証する。

20

【0032】

排他的暗号によれば一般に、平文名(ディレクトリエントリ中のファイル名またはディレクトリ名)が新しい名前にマップされる。マップされた名前は、任意選択で、ディケシファイ(decasified)された(大文字と小文字の区別をしない)名前と対応する大文字小文字情報にディケシファイされ、大文字と小文字の区別をしない複製名前検出が可能になる。マップ(および任意選択でディケシファイ)された名前は、次いで符号化され、暗号化される。この暗号化された名前(および任意選択で付随する大文字小文字情報)は、ディレクトリエントリの管理に責任を負うディレクトリグループに(例えば後に詳細に論じるパス名に基づいて)転送される。

30

【0033】

排他的暗号の詳細については、Douceur他の名義で2001年1月17日に出願され、Microsoft Corporationに譲渡された「Exclusive Encryption for a Secure Directory Service」という名称の同時係属米国特許出願第09/764962号を参照されたい。この出願は、参照によって本明細書に組み込まれる。

40

【0034】

(ファイルフォーマット)

図1のサーバレス分散ファイルシステム150のファイルフォーマットは、基本データストリームとメタデータストリームの2つの部分から成る。基本データストリームは、複数のブロックに分割されたファイルを含む。各ブロックは、対称暗号(例えばRC4)および暗号化鍵としてブロックのハッシュを使用して暗号化される。メタデータストリームは、ヘッダ、基本データストリーム中の暗号化されたブロックに索引を付けるための構造、およびいくつかのユーザ情報を含む。

【0035】

50

索引付けツリー構造は、各々のブロックのリーフノードを定義する。各リーフノードは、関連ブロックの復号に使用するアクセス値、および他のブロックとは独立に暗号化されたブロックを検証するのに使用する検証値から成る。一実施態様では、アクセス値は、ファイルブロックをハッシュし、得られたハッシュ値を、対称暗号およびランダムに生成された鍵を使用して暗号化することによって形成される。次いでこの鍵を、非対称暗号（例えば RSA）および暗号化鍵としてユーザの公開鍵を使用して暗号化する。検証値は、一方向ハッシュ関数（例えば SHA）を使用して、暗号化された関連ブロックをハッシュすることによって形成する。

【 0 0 3 6 】

ファイルのサイズに応じて、索引付け構造は、リーフノードをツリーブロックにグループ分けし、各ツリーブロックのハッシュ値を計算することによって形成された中間ノードを含むことができる。これらの中間ノードは、ブロックに再びセグメント化することができ、各ブロックをハッシュして次のノードを形成することができる。これを、ルートノードに到達するまで所望の回数繰り返すことができる。次いで、ルートノードをハッシュし、このハッシュ値を、メタデータヘッダおよびユーザ情報とともに使用して、ファイル全体の検証値を生み出す。一実施態様では、この全体ファイル検証値には、ユーザの署名が署名される。あるいは、このような署名なしでファイルを構築することもできる。

10

【 0 0 3 7 】

このファイルフォーマットは、ランダムに生成された鍵または任意のユーザ鍵についての知識なしで個々のファイルブロックを検証することをサポートする。ファイルのブロックを検証するため、ファイルシステムは任意選択で、全体ファイル検証値上の署名を評価し（存在する場合）、全体ファイル検証値がルートブロックのハッシュ、メタデータヘッダおよびユーザ情報と一致するかどうかをチェックし、次いで、検証する目的ブロックに関連した適当なリーフノードまでツリーを横断する。ファイルシステムが目的ブロックをハッシュし、ハッシュがリーフノードに含まれるアクセス値と一致する場合、そのブロックは本物である。

20

【 0 0 3 8 】

このファイルフォーマットはさらに、他のブロックを妨害することなく個々のブロックから読み取り、個々のブロックへ書き込むことをサポートする。このファイルフォーマットは、広大な非データ領域を有する粗なファイルに資する。

30

【 0 0 3 9 】

ファイルフォーマットの詳細については、Bolosky他の名義で2001年3月21日に出願され、Microsoft Corporationに譲渡された「On-Disk File Format for a Serverless Distributed File System」という名称の同時係属米国特許出願第09/814259号を参照されたい。この出願は、参照によって本明細書に組み込まれる。

【 0 0 4 0 】

（コンピューティング装置のアーキテクチャ）

図2に、図1の分散ファイルシステム150に参加しているコンピューティング装置102～106のうちの1台を表す例示的なコンピューティング装置200の論理的構成要素を示す。コンピューティング装置200は、サーバコンポーネント202、クライアントコンポーネント204、メモリ206、大容量記憶装置208および分散ファイルシステムインタフェース210を含む。典型的には、コンピューティング装置200は、追加の構成要素（例えばプロセッサ）を含むが、図面が乱雑にならないよう図2には、それらの追加構成要素を図示しない。さまざまなハードウェアおよびソフトウェアコンポーネントを有するより一般的なコンピュータアーキテクチャを後に図3を参照して説明する。

40

【 0 0 4 1 】

メモリ206は、RAM、ROM、フラッシュメモリなどの従来の任意の揮発性および/または不揮発性メモリとすることができる。大容量記憶装置208は、磁気ディスク、光ディスク、フラッシュメモリなどの従来の任意の不揮発性記憶装置とすることができる。大容量記憶装置208は、分散システム部分とローカル部分とに区画されている。図2に

50

は、大容量記憶装置 208 が 1 つしか示されていないが、コンピューティング装置 200 は、（異なるタイプまたは全て同じタイプの）複数の記憶装置 208 を含むことができる。

【0042】

コンピューティング装置 200 は、サーバレス分散ファイルシステム内で使用することを目的としたものであり、これに沿って、サーバコンポーネント 202 とクライアントコンポーネント 204 とを含む。サーバコンポーネント 202 は、記憶装置 208 に記憶された（または記憶する）ファイルまたはディレクトリエントリに関係する要求に装置 200 が応答しているときに、要求を処理し、クライアントコンポーネント 204 は、分散ファイルシステム内に記憶された（または記憶する）ファイルまたはディレクトリのために、装置 200 による要求の発行を処理する。クライアントコンポーネント 204 とサーバコンポーネント 202 は互いに独立に動作する。したがって、クライアントコンポーネント 204 が記憶中のファイルを、サーバコンポーネント 202 が大容量記憶装置 208 に記憶するといった状況も起こる可能性がある。

10

【0043】

クライアントコンポーネント 204 は、コンピューティング装置 150 に代って、ファイルおよびディレクトリの作成、記憶、検索、読取り、書込み、変更および検証目的のサーバレス分散ファイルシステム 150 へのアクセスを、インタフェース 210 とともに管理する記憶 / 検索制御モジュール 220 を含む。制御モジュール 220 は、ディレクトリグループ探索モジュール 222 を使用して、特定のファイルまたはディレクトリの管理に責任を負うディレクトリグループを識別し、ファイル暗号化モジュール 226 を使用してファイルを暗号化し、ディレクトリ暗号化モジュール 228 を使用して、ディレクトリエントリ中のファイル名およびディレクトリ名を暗号化する。これらのモジュールの動作については後に詳細に論じる。

20

【0044】

サーバコンポーネント 202 は、分散システム制御モジュール 250、複製識別子 252 およびサブツリー委任モジュール 254 を含む。分散システム制御モジュール 250 は、暗号化されたファイル 240 へのアクセスを管理する。分散システム制御モジュール 250 は、大容量記憶装置 208 と通信して、暗号化されたファイル 240 を記憶 / 検索する。分散システム制御モジュール 250 はさらに、コンピューティング装置 200（あるいはサーバレス分散ファイルシステム内の他の場所）に記憶されたディレクトリエントリ（図示せず）のレコードを、メモリ 206 および / または大容量記憶装置 208 内に維持する。サブツリー委任モジュール 254 は、サブツリーを他のディレクトリグループに委任するように動作する。これについては後に詳細に論じる。

30

【0045】

複製識別子 252 は、分散ファイルシステム内の暗号化された同一のファイルの識別を助ける。フォルトトレラント目的の意図的な複製ではない複製を見つけると、複製識別子 252 は、制御モジュール 250 に通知し、制御モジュール 250 は、この複製ファイルを排除し、排除されたファイルへのアクセス制御エントリを残りのファイルに追加する。

40

【0046】

図 3 に、分散ファイルシステムの実現に使用する、より一般的なコンピュータ環境 300 を示す。コンピュータ環境 300 は、コンピューティング環境の一例に過ぎず、コンピュータおよびネットワークアーキテクチャの使用または機能の範囲に関して何らかの制限を暗示しようとするものではない。コンポーネントの包含（または排除）、あるいは例示的なコンピュータ環境 300 内に示したコンポーネントの結合または組合せに関して、コンピュータ環境 300 が要件を有すると解釈してはならない。

【0047】

コンピュータ環境 300 は、コンピュータ 302 の形態の汎用コンピューティング装置を含む。コンピュータ 302 のコンポーネントには、1 または複数のプロセッサまたは処理ユニット 304、システムメモリ 306、およびシステムメモリ 306 にプロセッサ 30

50

4を含むさまざまなシステムコンポーネントを結合するシステムバス308が含まれる。ただしこれらに限定されるわけではない。

【0048】

システムバス308は、メモリバスまたはメモリコントローラ、周辺バス、アクセラレーテッドグラフィクスポート、およびさまざまなバスアーキテクチャを使用したプロセッサまたはローカルバスを含むいくつかのタイプのバス構造のうち、任意の1または複数のバス構造を表す。このようなアーキテクチャには例えば、Industry Standard Architecture (ISA)バス、Micro Channel Architecture (MCA)バス、Enhanced ISA (EISA)バス、Video Electronics Standards Association (VESA)ローカルバス、およびMezzanineバスとしても知られるPeripheral Component Interconnects (PCI)バスが含まれる。

10

【0049】

コンピュータ302は一般に、さまざまなコンピュータ読取り可能な媒体を含む。このような媒体は、コンピュータ302がアクセス可能な任意の使用可能媒体であり、揮発性および不揮発性媒体と、取外し可能および取外し不能媒体とを含む。

【0050】

システムメモリ306は、ランダムアクセスメモリ (RAM) 310などの揮発性メモリ、および/またはリードオンリーメモリ (ROM) 312などの不揮発性メモリの形態のコンピュータ読取り可能な媒体を含む。ROM 312には、スタートアップ時などに、コンピュータ302の要素間で情報を転送するのに助ける基本ルーチンを含む基本入出力システム (BIOS) 314が記憶されている。RAM 310は一般に、処理ユニット304が即座にアクセス可能なおよび/または現在操作中の、データおよび/またはプログラムモジュールを含む。

20

【0051】

コンピュータ302はさらに、取外し可能/取外し不能の揮発性/不揮発性コンピュータ記憶媒体を含むことができる。例えば図3には、取外し不能な不揮発性磁気媒体 (図示せず) の読取り/書込み用のハードディスクドライブ316、取外し可能な不揮発性磁気ディスク320 (例えば「フロッピー (登録商標) ディスク」) の読取り/書込み用の磁気ディスクドライブ318、およびCD-ROM、DVD-ROMなどの取外し可能な不揮発性光ディスク324の読取り/書込み用の光ディスクドライブ322が示されている。ハードディスクドライブ316、磁気ディスクドライブ318および光ディスクドライブ322は、1または複数のデータ媒体インタフェース360によってシステムバス308に接続されている。あるいは、ハードディスクドライブ316、磁気ディスクドライブ318および光ディスクドライブ322を、1または複数のインタフェース (図示せず) によってシステムバス308に接続することもできる。

30

【0052】

これらのディスクドライブおよびその関連コンピュータ読取り可能な媒体は、コンピュータ302のコンピュータ読取り可能な命令、データ構造、プログラムモジュールおよびその他のデータの揮発性記憶を提供する。この例には、ハードディスク316、取外し可能磁気ディスク320および取外し可能光ディスク324を示したが、コンピュータがアクセス可能なデータを記憶することができる、磁気カセットまたは他の磁気記憶装置、フラッシュメモリカード、CD-ROM、デジタルバーサタイルディスク (DVD) または他の光記憶装置、ランダムアクセスメモリ (RAM)、リードオンリーメモリ (ROM)、電氣的に消去可能なプログラマブルリードオンリーメモリ (EEPROM) などの他のタイプのコンピュータ読取り可能な媒体を利用しても、この例示的なコンピューティングシステムおよび環境を実現できることを理解されたい。

40

【0053】

ハードディスク316、磁気ディスク320、光ディスク324、ROM 312および/またはRAM 310上には、例えば、オペレーティングシステム326、1または複数のアプリケーションプログラム328、他のプログラムモジュール330およびプログラム

50

データ 3 3 2 を含む、任意の数のプログラムモジュールを記憶することができる。オペレーティングシステム 3 2 6、1 または複数のアプリケーションプログラム 3 2 8、他のプログラムモジュール 3 3 0 およびプログラムデータ 3 3 2 (またはこれらの組合せ) はそれぞれ、分散ファイルシステムをサポートする常駐コンポーネントの全部または一部を実現することができる。

【 0 0 5 4 】

ユーザは、キーボード 3 3 4、ポインティングデバイス 3 3 6 (例えば「マウス」) などの入力装置を介して、コマンドおよび情報をコンピュータ 3 0 2 に入力することができる。他の入力装置 3 3 8 (具体的には図示せず) には、マイクロホン、ジョイスティック、ゲームパッド、衛星アンテナ、シリアルポートおよび/またはスキャナなどが含まれる。これらの入力装置およびその他の入力装置は、システムバス 3 0 8 に結合された入出力インタフェース 3 4 0 を介して処理ユニット 3 0 4 に接続されるが、パラレルポート、ゲームポート、ユニバーサルシリアルバス (USB) などの他のインタフェースおよびバス構造によって接続することもできる。

10

【 0 0 5 5 】

モニタ 3 4 2 または他のタイプのディスプレイ装置を、ビデオアダプタ 3 4 4 などのインタフェースを介してシステムバス 3 0 8 に接続することもできる。モニタ 3 4 2 の他に、出力周辺装置として、入出力インタフェース 3 4 0 を介してコンピュータ 3 0 2 に接続することができるスピーカ (図示せず)、プリンタ 3 4 6 などのコンポーネントを含めることができる。

20

【 0 0 5 6 】

コンピュータ 3 0 2 は、リモートコンピューティング装置 3 4 8 などの 1 台または数台のリモートコンピュータへの論理接続を使用した、ネットワーク化された環境で動作することができる。リモートコンピューティング装置 3 4 8 は例えば、パーソナルコンピュータ、ポータブルコンピュータ、サーバ、ルータ、ネットワークコンピュータ、ピア装置または他の一般的なネットワークノードなどである。リモートコンピューティング装置 3 4 8 は、コンピュータ 3 0 2 に関して本明細書で説明した要素および特徴の多くまたは全てを含むことができるポータブルコンピュータとして示されている。

【 0 0 5 7 】

コンピュータ 3 0 2 とリモートコンピュータ 3 4 8 の間の論理接続は、ローカルエリアネットワーク (LAN) 3 5 0 および一般的なワイドエリアネットワーク (WAN) 3 5 2 として図示されている。このようなネットワーキング環境は、事業所、企業内コンピュータネットワーク、イントラネットおよびインターネットでよく見られる。

30

【 0 0 5 8 】

LAN ネットワーキング環境で実現するとき、コンピュータ 3 0 2 は、ネットワークインタフェースまたはネットワークアダプタ 3 5 4 を介してローカルネットワーク 3 5 0 に接続される。WAN ネットワーキング環境で実現するとき、コンピュータ 3 0 2 は一般に、ワイドネットワーク 3 5 2 を介して通信を確立するためのモデム 3 5 6 またはその他の手段を含む。モデム 3 5 6 は、コンピュータ 3 0 2 の内部または外部に設置することができる。入出力インタフェース 3 4 0 または他の適当な機構を介してシステムバス 3 0 8 に接続することができる。図示のネットワーク接続は例示的なものであり、コンピュータ 3 0 2 と 3 4 8 の間に通信リンクを確立する他の手段を使用することもできることを理解されたい。

40

【 0 0 5 9 】

コンピューティング環境 3 0 0 に示したようなネットワーク化された環境では、コンピュータ 3 0 2 に関して図示したプログラムモジュールまたはその一部分を、リモートメモリ記憶装置に記憶することができる。例えば、リモートコンピュータ 3 4 8 の記憶装置上に、リモートアプリケーションプログラム 3 5 8 が常駐する。図示の目的上、アプリケーションプログラムおよびオペレーティングシステムなどの他の実行可能プログラムコンポーネントを別個のブロックとして示したが、このようなプログラムおよびコンポーネントは

50

さまざまな時刻に、コンピューティング装置 302 のさまざまな記憶コンポーネントに存在し、コンピュータのデータプロセッサによって実行されることを理解されたい。

【0060】

分散ファイルシステム 150 の一実施態様は、1台または数台のコンピュータあるいはその他の装置によって実行されるプログラムモジュールなどのコンピュータ実行可能命令の一般的な文脈で記述することができる。プログラムモジュールは一般に、特定のタスクを実行し、または特定の抽象データ型を実現するルーチン、プログラム、オブジェクト、コンポーネント、データ構造などを含む。典型的には、さまざまな実施形態で、プログラムモジュールの機能を希望に応じて結合し、または分散させることができる。

【0061】

暗号化されたファイルのファイルフォーマットの一実施態様を、ある形態のコンピュータ読取り可能な媒体上に記憶し、またはある形態のコンピュータ読取り可能な媒体を横断して伝送することができる。コンピュータ読取り可能な媒体は、コンピュータがアクセス可能な任意の使用可能媒体とすることができる。コンピュータ読取り可能な媒体は例えば、「コンピュータ記憶媒体」および「通信媒体」を含む。ただしこれらに限定されるわけではない。

【0062】

「コンピュータ記憶媒体」は、コンピュータ読取り可能な命令、データ構造、プログラムモジュール、その他のデータなどの情報を記憶するため任意の方法または技術で実現された、揮発性および不揮発性の取外し可能および取外し不能媒体を含む。コンピュータ記憶媒体には、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、デジタルバーサタイルディスク(DVD)または他の光記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶装置または他の磁気記憶装置、あるいは所望の情報を記憶するのに使用することができ、コンピュータがアクセス可能な他の任意の媒体が含まれる。ただしこれらに限定されるわけではない。

【0063】

「通信媒体」は一般に、搬送波などの被変調データ信号または他の移送機構中に、コンピュータ読取り可能な命令、データ構造、プログラムモジュールまたは他のデータを具体化する。通信媒体にはさらに任意の情報送達媒体が含まれる。用語「被変調データ信号」は、信号中に情報を符号化するような方法で、その1または複数の特徴が設定または変更された信号を意味する。通信媒体には例えば、有線ネットワーク、直接有線接続などの有線媒体、音響、RF、赤外線などの無線媒体が含まれる。ただしこれらに限定されるわけではない。上記の任意の媒体の組合せもコンピュータ読取り可能な媒体の範囲に含まれる。

【0064】

(階層的記憶構造)

分散ファイルシステム 150 は、1または複数の名前空間ルートを含む階層的ファイル記憶構造を使用する。名前空間ルートはそれぞれ、1または複数のディレクトリまたはフォルダサブツリーをサポートすることができ、各サブツリーは、1または複数の追加のサブツリーをサポートすることができる。ディレクトリは、ゼロまたは1つ以上のファイル、および/またはゼロまたは1つ以上の他のディレクトリを保持することができる、シミュレートされたファイルフォルダと見ることができる。サブツリーは、1または複数のディレクトリを指し、ルートを含み(名前空間ルートを含むこともできる)、サブツリーのルートからサブツリーの全てのメンバへのパスがサブツリー自体の中にあるという特性を有する。図4に、ディレクトリA、B、C、D、E、F、G、H、J、I、M、K、Lを含む複数のサブツリーを有する名前空間ルートを含む、例示的な階層的名前空間400を示す。名前空間ルートのサブツリーには一般に、これよりも多くのディレクトリが含まれるが、説明しやすいように図4には少数のディレクトリだけを示した。

【0065】

各々のサブツリーは、1台または数台のコンピュータから成る、ディレクトリグループと呼ぶ一群のコンピュータによって管理される。本明細書では主に、サブツリーを管理する

10

20

30

40

50

ディレクトリグループとして論じるが、1または複数のディレクトリグループが、名前空間内の任意のディレクトリセットを管理することができる。図2の制御モジュール250など、コンピュータの1または複数のモジュールが、そのコンピュータが割り当てられたサブツリーを管理するディレクトリサービスを実現する責任を負う。一実施態様では、各ディレクトリグループが、後に詳細に論じるビザンチンフォルトトレラントグループ（単にビザンチングループとも呼ぶ）である。しかしディレクトリグループが、ビザンチンフォルトトレラントグループである必要はなく、他のグルーピングを使用することもできる。

【0066】

図4の実線は、ディレクトリ間の関係を示し、どのディレクトリがどのディレクトリのサブディレクトリであるかを識別する。例えば、ディレクトリCは、ディレクトリBのサブディレクトリである。あるディレクトリを、そのサブディレクトリの「親」ディレクトリと呼ぶことができる。例えば、ディレクトリBを、ディレクトリCの親ディレクトリと呼ぶことができる。

10

【0067】

図4の破線の囲みはそれぞれ、特定の破線の中に含まれるディレクトリを管理するディレクトリグループを示す。したがってこの例示的な名前空間400では、ルート名前空間がディレクトリグループ402によって管理され、ディレクトリA、B、C、FおよびGがディレクトリグループ404によって管理され、ディレクトリDおよびEがディレクトリグループ406によって管理され、ディレクトリHおよびJがディレクトリグループ408によって管理され、ディレクトリK、I、LおよびMがディレクトリグループ410によって管理される。

20

【0068】

特定のディレクトリまたは名前空間を管理するディレクトリグループは、そのディレクトリに記憶された各ファイルのディレクトリエントリ、ならびにそのディレクトリ内の各サブディレクトリのディレクトリエントリを維持する責任を負う。ファイルのディレクトリエントリの各々は、分散ファイルシステム150内の、そのファイルが記憶された1台または数台のコンピュータを識別する。サブディレクトリのディレクトリエントリの各々は、そのサブディレクトリを管理する責任を負ったディレクトリグループを識別する。ディレクトリエントリは、作成、変更およびアクセスのタイムスタンプ、読取り/書込みアクセス制御リスト、複製位置のセット、ファイルサイズなどの追加の情報を含むことができる。

30

【0069】

各ディレクトリグループは、名前空間ルートおよび/または名前空間内の1または複数のサブツリーを管理する責任を負う。各ディレクトリグループはさらに、1または複数の追加のサブツリーを識別し、それらの追加のサブツリーの管理責任を他のディレクトリグループに委任することができる。例えば、ディレクトリDおよびEは元々、ディレクトリグループ404によって管理されていたが、後に、ディレクトリグループ406に委任された。

【0070】

ディレクトリグループはいつでも、サブツリーを別のディレクトリグループに委任することを決定することができる。一実施態様では、この決定が作業負荷に基づき、ディレクトリグループは、自身が過負荷になりつつあると判定したときにサブツリーの委任を決定する。ディレクトリグループは、自身が過負荷になりつつあることをさまざまな因子を使用して判定することができ、例示的な一実施態様では、各ディレクトリグループは、予想されるマシン1台あたりの平均ディレクトリ数（例えば10,000程度）にほぼ等しいサイズのサブツリーを管理しようとする。

40

【0071】

サブツリーを委任する先のディレクトリグループは、さまざまな方法で決定することができる。一実施態様では、委任を実行する側のディレクトリグループが、分散ファイルシス

50

テム150内の知っているコンピュータの中からランダムに選択し、選択したコンピュータを、サブツリーを委任する先の新しいディレクトリグループとして使用する。他のさまざまな因子をこの選択プロセスに影響させることができる（例えば、可用性の低いコンピュータは選択しない、最近にサブツリーを委任したコンピュータは選択しないなど）。

【0072】

ディレクトリグループは、1または複数のメンバによってデジタル署名された委任証明書を生成することによって、特定のサブツリーを委任することができる。複数のメンバが委任証明書に署名する状況では、署名プロセスをさまざまな形式で実施することができる。一実施態様では、各メンバは、委任証明書のコピーに署名する。他の実施態様では、委任証明書に繰り返し署名される（例えば、1つのメンバが証明書に署名し、次いで、デジタル署名されたこの証明書に他のメンバが署名する）。デジタル署名を検証するときには検証者がその順番を知っていさえすれば、異なるメンバが証明書に署名する順番は重要ではない（例えば、署名順を用いて検証者を前もってプログラムしておくこと、または順番を識別する情報を証明書に含めることができる）。以下に、4つの署名者によって繰り返し署名された例示的な証明書を示す。

$s_4(s_3(s_2(s_1(DC))))$

ここで、DCは、デジタル署名される委任証明書を表し、 $s_i()$ は、署名者iが()の内容にデジタル署名したことを指示する。

【0073】

一実施態様では、ディレクトリグループのメンバ（コンピュータ）の数が、設計者が許容できるようにしたいと考える障害コンピュータの数によって決まる。本明細書では障害コンピュータが、アクセス不能のコンピュータ（例えば電源が切られた、または誤動作しているコンピュータ）、または不正使用されたコンピュータ（例えば、悪意あるユーザまたはプログラムがコンピュータにアクセスし、適当な応答を与えるまたは不適切なデータを与えるなど、問合せに対して不適切に応答することができるコンピュータ）を指す。具体的な1つの例では、f台の障害コンピュータを許容するために、ディレクトリグループは、 $3f + 1$ 台のコンピュータを含む。さらにこの例では、少なくとも $f + 1$ 台のコンピュータが委任証明書にデジタル署名する。

【0074】

各々の名前空間ルートには、認証局(CA)から取得した証明書が関連づけられている。認証局は、名前空間の作成を検証する信用機関である。サブツリーに関連づけられた各委任証明書は、カレントサブツリーから、ゼロまたは1つ以上の他のサブツリーを介してCAが署名した名前空間ルート証明書までさかのぼる証明書のチェーンを含む。したがって、各委任証明書には、そのサブツリーを管理することを許可されたディレクトリグループであることを(CAが署名した証明書までさかのぼる証明書のチェーンを確立することによって)証明する複数の証明書が関連づけられている。

【0075】

委任証明書は、さまざまな構成要素を含むことができ、一実施態様では委任証明書が、(1)委任を実行するディレクトリグループが管理しているサブツリーのルートよりも下にある、委任するパスの識別、(2)委任を実行するディレクトリグループに委任されたサブツリーのルートの識別、(3)委任するサブツリーの識別、および(4)サブツリーを委任する先のグループのメンバの識別を含む。サブツリーおよびパスメンバの識別はさまざまであり、実際のディレクトリ名（例えばディレクトリA、B、C、Dなどの名前）、あるいは識別番号（例えば、Globally Unique Identifier (GUID)）とすることができる。識別番号を使用すると、ディレクトリ名が変更された場合に、委任証明書を作成し直す必要性を回避することができる。

【0076】

図4を参照して委任証明書の例を見つめる。ディレクトリグループ402は、グループ402が名前空間ルートを管理する権限を有することを証明する証明書をCAから取得する。この証明書は以下の形式をとる。

10

20

30

40

50

$ourcA (Root, GUID_{Root}, DG_{402})$ (1)
 ここで、 $ourcA$ は、この証明書がCA「OurCA」によって署名されたことを指示し、 $Root$ は名前空間ルートの名前、 $GUID_{Root}$ は名前空間ルートのGlobally Unique Identifierであり、 DG_{402} は、ディレクトリグループ402のメンバの名前（または他の識別子）を表す。

【0077】

ディレクトリグループ402がディレクトリAから始まるサブツリーをディレクトリグループ404に委任すると決定すると、ディレクトリグループ402は、ディレクトリグループ404のメンバに渡す委任証明書を生成する。この委任証明書は、先の証明書(1)ならびに以下の証明書を含む。

$DG_{402} (GUID_{Root}/A, GUID_A, DG_{404})$ (2)
 ここで、 DG_{402} は、ディレクトリグループ402のメンバがこの証明書に署名したことを指示し、 $GUID_{Root}/A$ は、ディレクトリグループ402に委任されたサブツリーのルートのGUID($GUID_{Root}$)とディレクトリグループ404に委任されるパス(/A)、 $GUID_A$ は、委任されるサブツリー（すなわちディレクトリAから始まるサブツリー）のGlobally Unique Identifierであり、 DG_{404} は、ディレクトリグループ404のメンバの名前（または他の識別子）を表す。

【0078】

同様に、ディレクトリグループ404がディレクトリDから始まるサブツリーをディレクトリグループ406に委任すると決定すると、ディレクトリグループ404は、ディレクトリグループ406のメンバに渡す委任証明書を生成する。この委任証明書は、先の証明書(1)および(2)、ならびに以下の証明書を含む。

$DG_{404} (GUID_A/B/C/D, GUID_D, DG_{406})$ (3)
 ここで、 DG_{404} は、ディレクトリグループ404のメンバがこの証明書に署名したことを指示し、 $GUID_A/B/C/D$ は、ディレクトリグループ404に委任されたサブツリーのルートのGUID($GUID_A$)とディレクトリグループ406に委任されるパス(/B/C/D)、 $GUID_D$ は、委任されるサブツリー（すなわちディレクトリDから始まるサブツリー）のGlobally Unique Identifierであり、 DG_{406} は、ディレクトリグループ406のメンバの名前（または他の識別子）を表す。

【0079】

図示の例では、委任証明書が、特定のサブツリー内の各ディレクトリに対して発行されるのではなく、委任ポイントで発行される。例えば、委任証明書は、A（サブツリーの最上位ディレクトリ）に対して発行されるのであって、/A/Bまたは/A/B/Cに対して発行されるのではない。

【0080】

図5は、サブツリーの管理責任を他のディレクトリグループに委任する例示的なプロセス500を示す流れ図である。プロセス500は、サブツリーの管理責任の委任を実行しているディレクトリグループ内のコンピュータのサブツリー委任モジュール254によって実行される。最初に、サブツリーを委任する先の一群のコンピュータを識別する（アクト502）。そのサブツリーに対する委任証明書を生成し（アクト504）、委任するグループの1または複数のメンバがデジタル署名する（アクト506）。次いで、そのサブツリーの管理責任を委任する先のコンピュータ群に、デジタル署名された委任証明書を発行する（アクト508）。

【0081】

図4に戻る。分散ファイルシステム150内の各コンピュータは、名前空間内のパス名のあるサブセットを、そのパス名を管理するディレクトリグループにマップするローカルキャッシュ（例えば図2のキャッシュ260）を維持する。例えば、ある特定のコンピュータのキャッシュが、パス名/A、/A/B、/A/B/C、/A/Fおよび/A/F/Gを、ディレクトリグループ404にマップするマッピングを含む。コンピュータは、それらのキャッシュの中にさまざまなマッピングを有することができるが、典型的には、ディ

10

20

30

40

50

レトリグループ（ディレクトリグループ402）を管理する名前空間ルートに、名前空間ルートをマッピングすることを少なくとも含む。

【0082】

パス名から管理ディレクトリグループへのマッピングを維持することによって、コンピュータが、ディレクトリグループ探索プロセスの少なくとも一部分を、名前空間ルートを管理しているディレクトリグループ（およびおそらく他のディレクトリグループ）にアクセスする必要なしに、それ自体でローカルに実行することができる。例えば、あるコンピュータが、パス名/A/B/foo.txtを有する「foo.txt」と呼ばれるファイルにアクセスすることを希望しており、コンピュータがそのローカルキャッシュに、ディレクトリグループ404に対するパス名のマッピングを有すると仮定する。この例では、コンピュータが、ディレクトリBの中のファイルを管理し、したがってファイルfoo.txtを管理するディレクトリグループ404内のメンバを、自身のローカルキャッシュから容易に識別することができる。したがって、コンピュータは、ファイル「foo.txt」の位置を決定するのにどのコンピュータにアクセスすべきか（すなわち、パス名/A/Bのディレクトリエントリをどのコンピュータが管理しているか）を、ディレクトリグループ402または404にアクセスする必要なしに、キャッシュの中の情報に基づいて決定する。

10

【0083】

コンピュータがディレクトリグループへのパス名全体をマップするだけの十分な情報を、ローカルキャッシュ中に有していない場合、コンピュータは、そのキャッシュに存在するそのパス名の最も長いプリフィックスのマッピングを見つける。次いでコンピュータは、最も長いプリフィックスの最後のディレクトリを管理するディレクトリグループにアクセスして、パス名の残りの部分および委任証明書を、できるだけ多く管理するディレクトリグループを決定する。ディレクトリグループにアクセスし、委任証明書を取得するこのプロセスは、適当なマッピングが見つかるまで続けられる。

20

【0084】

例えば、あるコンピュータが、パス名/A/B/C/D/foo2.txtを有する「foo2.txt」と呼ばれるファイルにアクセスすることを希望しており、コンピュータがそのローカルキャッシュに、ディレクトリグループ404に対するパス名のマッピングは持っているが、ディレクトリグループ406に対しては持っていないと仮定する。コンピュータはパス名を見て、そのキャッシュの中にあるパス名（/A/B/C）における最も長いプリフィックスに対するマッピングを見つけ、ディレクトリの管理に責任を負うディレクトリグループ、すなわちディレクトリグループ404にアクセスする。コンピュータは、ディレクトリグループ404のメンバに、パス名/A/B/C/D/foo2.txtの関連サブツリーの委任証明書、すなわちディレクトリグループ406の委任証明書を問い合わせる。ディレクトリグループ404のメンバは、問い合わせたコンピュータにこの委任証明書を戻し、このコンピュータは、委任証明書を（例えば署名したコンピュータの公開鍵に基づいて）検証することができる。受け取られた委任証明書は、ディレクトリ/Dの管理に責任を負うディレクトリグループを識別し、そのためコンピュータは、ファイル「foo2.txt」がある場所を決定するために、そのディレクトリグループにアクセスすればよいことを知る。このように、ファイル「foo2.txt」の位置を決定するためにどのコンピュータにアクセスすべきかを決定することには、ディレクトリグループ404のメンバにアクセスすることが含まれるが、この決定のためにディレクトリグループ402のメンバへアクセスする必要はない。

30

40

【0085】

図6は、特定のパス名の管理に責任を負うディレクトリグループを探索する例示的なプロセス600を示す流れ図である。プロセス600は、探索中のパス名にアクセスすることを希望するコンピュータのディレクトリグループ探索モジュール222（図2）によって実行される。最初に、ディレクトリグループへのマッピングのローカルキャッシュにアクセスし（アクト602）、このキャッシュにパス名の最も長いプリフィックスに対するマ

50

ッピングを見つけ出す(アクト604)。次いで、パス名全体がマップされているかどうかに基づいて処理を続行する(アクト606)。パス名全体がマップされている場合、ディレクトリグループ探索プロセスは完了となる(アクト608)。しかし、パス名全体がマップされていない場合には、パス名の最後にマップされたプリフィックスを管理するグループのメンバから、関連サブツリーの委任証明書を取得する(アクト610)。次いで、受け取った委任証明書を検証する(アクト612)。委任証明書が正しく検証されないか、または委任証明書を取得することができない場合には、アクト610へ戻り、グループの別のメンバを選択して問い合わせる。最後にマップされたプリフィックスを管理しているグループの少なくとも1つのメンバが正常に機能している限り、このプロセスは最後には成功する。最後にマップされたプリフィックスを管理しているグループのメンバが1つ

10

も正常に機能していない場合、プロセスは、有効なより短いプリフィックスを探るか、または名前空間ルートへ戻る。委任証明書のチェーンが検証されたら、証明書からのパス名マッピング情報をローカルキャッシュに追加する(アクト614)。次いで、プロセスはアクト606へ戻る。このときの最も長いプリフィックスは、以前の最も長いプリフィックスに新しい関連サブツリー情報を連結したものになる(例えば、パス名が/A/B/C/D/E/F、以前の最も長いプリフィックスが/A/B、新しい関連サブツリーが/C/Dである場合、新しい最も長いプリフィックスは/A/B/C/Dとなる)。次いで、パス名全体がマップされるまで、アクト606、610、612および614を繰り返す。異なるディレクトリの管理を異なるディレクトリグループ上へ分離することによって、管理責任は、異なる複数のコンピュータに分散する。これによって、特定のコンピュータ

20

、特に名前空間ルートのディレクトリグループおよび名前空間ルートに最も近いディレクトリグループのコンピュータにかかる管理負担は低減する。例えば、特定のパス名をルートノードから解析する必要はなく、そのローカルキャッシュを介して、そのパス名の途中までは拾い上げることができる。

【0086】

(ディレクトリおよびファイルの複製および記憶)

図1の分散ファイルシステム150は、ディレクトリエントリの記憶とディレクトリエントリに対応するファイルの記憶を異なる方法で管理する。システム150内に記憶するファイルは複製され、システム150内の異なる複数のコンピュータに保存される。さらに、そのファイルに対するディレクトリエントリが生成され、ピザンチンフォルトトレラントグループの一部であるシステム150内の異なる複数のコンピュータに保存される。ディレクトリエントリは、後に詳細に論じるように、ファイルが保存されるよりも多くのコンピュータに保存される。

30

【0087】

ファイルおよびディレクトリエントリを記憶する本明細書に記載のさまざまな処理を、先に論じた階層的記憶構造とともに使用することができる。しかし、ファイルおよびディレクトリエントリを記憶する本明細書に記載のさまざまな処理を、階層的記憶構造を使用しないシステムで使用することもできる。

【0088】

ピザンチンフォルトトレラントグループは、たとえそのうちのいくつかのコンピュータに障害(故障または他の理由で使用できない)があっても、情報を記憶し、かつ/または他のアクションを実行するのに使用することができる一群のコンピュータである。コンピュータは、コンピュータを操作する悪意のユーザ、コンピュータ上で走る悪意のプログラムなど、さまざまな方法で危険にさらされる。これらのコンピュータからは、リクエストへの応答を拒否する、リクエストに対して不正確なまたは不要な情報で故意に応答するなど、任意のタイプの振舞いを観察することができる。このようなコンピュータの存在にもかかわらず、ピザンチンフォルトトレラントグループは、情報を正確に記憶し、かつ/または他のアクションを実行することができる。ピザンチングループは当業者に周知であり、したがって、本発明に関する場合を除き、さらに論じることはしない。

40

【0089】

当業者には周知のとおり、あるタイプの計算では、 f 台の故障コンピュータの存在にもかかわらず（故障コンピュータは危険にさらされ、またはパワーダウンなどの他の理由で使用不能である）、正確に動作できるようにするためには、ビザンチンフォルトトレラントグループが、少なくとも $3f + 1$ 台のコンピュータを含んでいなければならない。分散ファイルシステム 150 では、ディレクトリエントリが、ビザンチンフォルトトレラントグループの $3f + 1$ 台のコンピュータに記憶され、ファイル自体は、 $f + 1$ 台のコンピュータ（ディレクトリエントリが記憶された同じコンピュータのうちの 1 台または数台とすることができる）に記憶される。

【0090】

図 7 に、ファイルおよび対応するディレクトリエントリをサーバレス分散ファイルシステム内に記憶する例示的な記憶を示す。ファイルシステム 700（例えば図 1 のサーバレス分散ファイルシステム 150）は、12 台のコンピュータ 702、704、706、708、710、712、714、716、718、720、722 および 724 を含む。システム 700 の設計者が、2 台のコンピュータの故障を許容できるようにしたいと考えているとすると、ビザンチンフォルトトレラントグループは、少なくとも 7 台（ $(3 \times 2) + 1$ ）のコンピュータを含まなければならない。コンピュータ 702 ~ 714 を含むビザンチングループ 726 が示されている。

【0091】

ファイル 728 をファイルシステム 700 に記憶するときには、対応するディレクトリエントリ 730 が、適当なディレクトリグループ（ファイル 728 のパス名に基づいてファイルが記憶されたディレクトリの管理に責任を負うディレクトリグループ）内のコンピュータによって記憶される。ディレクトリエントリ 730 に対する図 7 のディレクトリグループは、ビザンチングループ 726 であり、そのため、ディレクトリエントリ 730 は、ビザンチングループ 726 の正常に機能している各々のコンピュータ 702 ~ 714 上に記憶される。したがってディレクトリエントリ 730 は、最大 7 台の異なるコンピュータ上に記憶される。一方、ファイル 728 は複製され、3 台のコンピュータ（コンピュータ 716、720 および 724）にそれぞれ記憶される。図示のとおり、ファイル 728 を記憶したコンピュータが、ビザンチングループ 726 のコンピュータである必要はなく、一般に、ビザンチングループ 726 のコンピュータではない（ただし、任意選択で、ファイル 728 が記憶された 1 台または数台のコンピュータがビザンチングループ 726 のコンピュータであってもよい）。

【0092】

各ディレクトリエントリは、対応するファイルの名前、ファイルが記憶されたコンピュータの識別、およびファイルの内容がディレクトリエントリに対応するかどうかを検証することを可能にするファイル検証データを含む。ファイル検証データは、異なるさまざまな形式をとることができ、一実施態様では、ファイル検証データが、MD5（Message Digest5）、SHA-1（Secure Hash Algorithm-1）などの暗号法上安全なハッシュ関数をファイルに適用することによって生成されたハッシュ値である。記憶装置からファイルを検索するとき、検索するコンピュータは、ハッシュ値を再生成し、その値をディレクトリエントリ中のハッシュ値と比較して、コンピュータが正確なファイルを受け取ったことを検証することができる。他の実施態様では、ファイル検証データが、ファイル識別番号（例えばファイルの一意識別子）、ファイルのバージョン番号、およびファイル上に署名があるユーザの名前を結合したものである。

【0093】

図 8 は、サーバレス分散ファイルシステムにファイルを記憶する例示的なプロセスを示す流れ図である。最初に、クライアントコンピューティング装置が新しいファイル記憶要求を受け取る（アクト 802）。クライアントは、ファイルおよびファイル名を暗号化し、ファイル内容のハッシュを生成する（アクト 804）。クライアントは、暗号化されたファイル名およびファイル内容ハッシュを、ディレクトリエントリを作成するよう求める要求とともに、適当なビザンチンフォルトトレラントディレクトリグループに送る（アクト

10

20

30

40

50

806)。ディレクトリグループは、要求の有効性を、例えば、ファイル名が既存の名前と競合しないこと、およびクライアントが要求の内容を実行する許可を有することを検証することによって、調べる(アクト808)。要求が有効と認められなかった場合、要求は失敗に終わる(アクト810)。しかし、要求が有効と認められた場合には、ディレクトリグループが、この新しいファイルのディレクトリエントリを生成する(アクト812)。ディレクトリグループはさらに、この新しいファイルの複製セットを決定し、その複製セットを、新しく生成されたディレクトリエントリに追加する(アクト814)。ファイルの複製も生成し(アクト816)、ファイルシステム内の複数のコンピュータに保存する(アクト818)。

【0094】

ビザンチングループにディレクトリエントリを記憶し、このエントリにファイル検証データを含めることによって、フォルトトレランス(最大f故障まで)が維持される。しかし、ファイルをディレクトリとは別に記憶し、ファイルへのアクセスにビザンチン操作を使用しないことによって、記憶空間の必要量およびビザンチン操作が低減される。例えば、ディレクトリエントリは、100バイト程度であるのに対して、ファイル自体は、数千または数十億バイトにもなる。

【0095】

(ディレクトリおよびファイルのロック機構)

図1の分散ファイルシステム150の各オブジェクト(例えばディレクトリおよびファイル)は、一組の専用ロックに関連づけられている。これらのロックは、アプリケーションが実行したいオペレーションのタイプに基づいて、そのオペレーションを実行するためにアプリケーションがディレクトリまたはファイルを開くことができるか否かを決定するのに使用される。ロックは、ロックのタイプおよび競合のレベルによって決まる特定の時間範囲を有するリース(lease)と見ることができる。例えば、書込みロックの時間範囲は数分であり、一方、読取りロックの時間範囲は数日にもなる。アプリケーションが、オブジェクトに対してオペレーションを実行したいとき、このアプリケーションを実行中のクライアントコンピュータは、そのオペレーションを実行するのに必要なロックをすでに有しているかどうかを調べる。必要なロックを持っていない場合には、そのオブジェクトの管理に責任を負うディレクトリグループから適当なロックを要求する。所望のオペレーションを実行し終わると、アプリケーションは任意選択で、取得したロックを解放し、あるいはロックが自動的に期限切れになるか、または管理ディレクトリグループによってロックがリコールされるまで、ロックを保持する。

【0096】

特定のディレクトリに対して、そのディレクトリを実装するビザンチンフォルトトレラントグループが、そのディレクトリ内の全てのファイル、そのディレクトリのサブディレクトリの名前、およびそのディレクトリ自体を削除する権利に対するロックを制御する。ロック機構は、適当なファイルおよびディレクトリの幅広い(細分性の粗い)ロックを、要求しているクライアントコンピュータに許可しようとし、クライアントコンピュータは、ロックを取得するのに複数のビザンチンメッセージを要求するのではなく、1回のビザンチンロックの取得で多くの読取りおよび/または更新を処理することができる。

【0097】

図示の例では、ロック機構が10個の異なるロック、すなわち読取り、書込みク、オープン読取り(Open Read)、オープン書込み(Open Write)、オープン削除(Open Delete)、非共用読取り(Not Shared Read)、非共用書込み(Not Shared Write)、非共用削除(Not Shared Delete)、挿入、および排他ロックを使用する。読取りおよび書込みロックは、オブジェクト中のデータ(例えばファイルの内容)へのアクセスを制御するのに使用される。オープン読取り、オープン書込み、オープン削除、非共用読取り、非共用書込み、および非共用削除ロックは、オブジェクトのオープンを制御するのに使用される。挿入および排他ロックは特殊用途のロックである。これらの10個のロックについては後にさらに詳細に論じる。アプリケーションは、実行したいオペレーションに応じて、これ

10

20

30

40

50

らのロックの中から適当なものを要求する。

【0098】

[読取りロック] 読取りロックは、アプリケーションが関連ファイルを読み取ることができるようアプリケーションが要求する。読取りロックは、書込みロックとともに、ディレクトリグループがファイル中のデータを矛盾なしに保つことを可能にする。

【0099】

[書込みロック] 書込みロックは、アプリケーションが関連ファイルに書き込む（ファイルを更新するとも言う）ことができるようアプリケーションが要求する。書込みロックは、読取りロックとともに、ディレクトリグループがファイル中のデータを矛盾なしに保つことを可能にする。

10

【0100】

アプリケーションがオブジェクトを開きたいとき、ディレクトリグループは、2つのチェックを実行する。(1)アプリケーションが求めているモードが、すでにそのオブジェクトを開いている他のアプリケーションと競合するかどうか、および(2)アプリケーションがオブジェクトを共用する気があるオペレーションが、他のアプリケーションがそのオブジェクトをすでに開き、そのオブジェクトを共用する気があることを指示したオペレーションと競合するかどうかをチェックする。10個のロックのうち、オープン読取り、オープン書込み、オープン削除、オープン非共用読取り、オープン非共用書込み、およびオープン非共用削除の6つのロックが、このチェックをサポートすることに向けられる。これらのロックは、オブジェクトを開く能力をアプリケーションに許可するのに使用されるが、そのオブジェクトのデータを取得できることを保証するとは限らない（データにアクセスするためには、（アプリケーションが実行したいオペレーションのタイプに応じて）読取りロックまたは書込みロックを取得する）。

20

【0101】

[オープン読取りロック] オープン読取りロックは、読取りのためアプリケーションが関連オブジェクトを開くことができるよう、アプリケーションが要求する。

【0102】

[オープン書込みロック] オープン書込みロックは、書込みのためアプリケーションが関連オブジェクトを開くことができるよう、アプリケーションが要求する。

【0103】

[オープン削除ロック] オープン削除ロックは、削除のため、アプリケーションが関連オブジェクトを開くことができるよう、アプリケーションが要求する。

30

【0104】

[オープン非共用読取りロック] オープン非共用読取りロックは、アプリケーションがオブジェクトを読み取る能力を他のアプリケーションと共用したくないときにアプリケーションが要求する。

【0105】

[オープン非共用書込みロック] オープン非共用書込みロックは、アプリケーションがオブジェクトに書き込む能力を他のアプリケーションと共用したくないときにアプリケーションが要求する。

40

【0106】

[オープン非共用削除ロック] オープン非共用削除ロックは、アプリケーションがオブジェクトを削除する能力を他のアプリケーションと共用したくないときにアプリケーションが要求する。

【0107】

サポートされる他の2つのロックが挿入ロックおよび排他ロックである。

【0108】

[挿入ロック] 挿入ロックは、ディレクトリ内のオブジェクトに対して特定の名前を作成するためにアプリケーションが要求する。挿入ロックの許可は、特定の名前を有するオブジェクトを作成する許可をアプリケーションに与える。挿入ロックは、同じオブジェクト

50

名を有する他の挿入ロック、およびディレクトリの排他ロックと競合する。

【0109】

[排他ロック] 排他ロックは、上述した9つのロック全てを取得するためにアプリケーションが要求し、ディレクトリ内に存在することができる(しかしまだ存在しない)各々の可能な名前の挿入ロックを含む。ディレクトリの排他ロックは、そのディレクトリのファイルまたはサブディレクトリの排他ロックを意味せず、ディレクトリの名前空間だけを意味する。排他ロックは、先に論じた9つのロックのそれぞれと競合する。

【0110】

異なるさまざまなロック間には、さまざまな競合が存在する。表1は、例示的な一実施態様におけるロック間の競合を示す競合マトリックスである。表1では以下の省略形を使用した: Ins (挿入)、Excl (排他)、O-R (オープン読取り)、O-W (オープン書込み)、O-D (オープン削除)、O-!R (オープン非共用読取り)、O-!W (オープン非共用書込み) および O-!D (オープン非共用削除)。表1の欄の中の「X」は、対応する2つのロック間の競合を示す。例えば、オープン読取りは、オープン非共用読取りと競合するが、オープン非共用書込みとは競合しない。

10

【0111】

【表1】

表1

	Ins	Read	Write	Excl	O-R	O-W	O-D	O-!R	O-!W	O-!D
Ins	X	X	X	X						
Read	X		X	X						
Write	X	X	X	X						
Excl	X	X	X	X	X	X	X	X	X	X
O-R				X				X		
O-W				X					X	
O-D				X						X
O-!R				X	X					
O-!W				X		X				
O-!D				X			X			

20

30

【0112】

図9は、特定のオブジェクトを開くことを許可するかどうか決定する例示的なプロセスを示す流れ図である。図9のプロセスは、その特定のオブジェクトの管理に責任を負うディレクトリグループによって実現される。図9のプロセスでは、特定のオブジェクトを開くことを要求しているクライアントが、所望のオブジェクトを開くのに必要なロックをまだ持っていないと仮定する。最初に、識別された特定のロックを有するオブジェクトにアクセスするリクエストを受け取る(アクト902)。選択されたロックが意味するモードが、別のクライアントに許可されたロックと競合するかどうかに関するチェックが、ディレクトリグループによって実施される(アクト904)。例えば、このリクエストが、読取りのためオブジェクトを開くリクエストであり、他のアプリケーションが、非共用読取りロックを有するこのオブジェクトをすでに開いている場合、選択されたロックが意味するモード(オープン読取り)は、そのオブジェクトをすでに開いている他のアプリケーションと競合する。ディレクトリグループは、競合するロックをクライアントに発行したかどうかを知っているだけで、そのクライアントが現在、オブジェクトへのアプリケーションのアクセスを可能にするためにそのロックを使用しているかどうかを知らないの、いくつかのケースでは、アクト904のチェックを実施するのに、ロックを現在保持している

40

50

クライアントがそれを放棄する気があるかどうかを尋ねる必要がある。

【0113】

アクト904のチェックで競合が識別されなかった場合、要求されたロックが許可され、選択したロックを有するファイルをアプリケーションが開くことができるようになり(アクト906)、アクト902のリクエストが許可される。次いで、これらのロックが許可されたこと、およびそれら許可されたクライアントが、ディレクトリグループによって保存され(アクト908)、後続のリクエストに対して競合を決定し、必要に応じてロックのリコールを試みることができる。

【0114】

しかし、アクト904のチェックで競合が識別された場合には、競合するロックを保持しているクライアントに、それらを戻すよう求めるリクエストが発行される(アクト910)。次いで、要求されたロックが全て戻されたかどうかに関するチェックが実施される(アクト912)。要求されたロックが全て戻された場合には、要求のロックが許可され、選択したロックを有するファイルをアプリケーションが開くことができるようになり(アクト906)、これらのロックが記録される(アクト908)。一方、要求されたロックのうち戻されないロックがある場合には、ディレクトリグループは、このオープンリクエストを拒絶する(アクト914)。

【0115】

1台だけのクライアントコンピュータが名前空間のある領域にアクセスするとき性能を向上させる試みにおいて、ファイルシステム150は、アプリケーション(またはクライアント)が近い将来に、追加の関連ロックを要求する可能性が高いとの前提のもとに、クライアント上で実行中のアプリケーションが要求するよりも広い範囲のロックを発行することができる。例えば、アプリケーションがファイル/A/B/C/foo.txtを開く場合、クライアントは、このファイルに対するロックを要求する。ディレクトリグループがそのロックを許可する場合、ディレクトリグループはそのロックを、/A/B/Cのディレクトリロックにアップグレードすることができる(例えば、過去の実施に基づいて、ディレクトリ上の競合がまれであるとディレクトリグループが判定した場合)。次いでアプリケーションが、同じディレクトリ内の別のファイルを開く場合、クライアントは、ディレクトリグループから他のロックを要求する必要なしにそのファイルを開くことができる。

【0116】

クライアントのロックリクエストが、他のクライアントに許可された既存のロックと競合した場合、ディレクトリグループは、(例えば、アクト914でリクエストを拒絶するのではなく)アクト910で、新しいリクエストと競合しないロックに、以前に発行されたロックをダウングレードすることを試みることができる。ロックのアップグレードの結果、クライアントは、要求していないロックを保持するので、ロックのダウングレードが成功する可能性は一般に低くない。このロックリコールに失敗した場合、リクエストは拒絶される。

【0117】

ファイルシステム内のオブジェクトにさまざまなオペレーションを実行することができる。表2に、より一般的ないくつかのオペレーション、およびオペレーションを実行するためにアプリケーションが要求するロックを記載する。

【0118】

【表2】

10

20

30

40

表2

オペレーション	説明
オブジェクト読取り	ディレクトリまたはファイルを読み取るためのリクエスト。オブジェクトのオープン読取りロックおよびそれに続く読取りロックを必要とする。アプリケーションは任意選択で、任意のオープン非共有ロックを要求することができる。
オブジェクト書込み ／更新	ファイルに書き込むためのリクエスト。オブジェクトのオープン書込みロックおよびそれに続く書込みロックを必要とする。要求されたならば、アプリケーションは任意選択で、任意のオープン非共有ロックを要求することができる。
ファイル削除	ファイルを削除するためのリクエスト。オープン削除および書込みロックを必要とする。通常、アプリケーションはさらに、全てのオープン非共有ロックを要求する。
ディレクトリ削除	ディレクトリを削除するためのリクエスト。ディレクトリに対する排他ロックを必要とする。ディレクトリを削除できるのはディレクトリが空のときだけである。
ディレクトリ名変更	ディレクトリ名を変更するためのリクエスト。親ディレクトリの排他ロック（名前を変えられているディレクトリ名が変更されるディレクトリはサブディレクトリである）、およびあて先ディレクトリ中の新しいディレクトリ名に対する挿入ロックを必要とする。名前変更が複数のディレクトリにまたがる場合、新しい親ディレクトリに対する挿入ロックが必要となる。
ファイル名変更	ディレクトリ中のファイルの名前を変更するためのリクエスト。ファイルの書込みロック、およびディレクトリ（名前変更が複数のディレクトリにまたがる場合には別のディレクトリとすることができる）中の新しい名前に対する挿入ロックを必要とする。
オブジェクト作成	新しいファイルまたはディレクトリを作成するためのリクエスト。新しい名前に対する挿入ロックを必要とする。

【0119】

ファイルへの変更は、そのコンピュータによってローカルに実施され、次いでファイルは（暗号化された後には）、そのファイルの管理に対して責任を負うディレクトリグループにプッシュバックされる。この情報は、ディレクトリグループ内のさまざまなコンピュータに記憶され、更新されたファイルは、適当なコンピュータに記憶される。

【0120】

（結語）

以上の説明では、構造上の特徴および／または方法論的行為に特有の言葉を使用したが、添付の請求項に定義した発明は、記載した特定の特徴または行為に限定されないことを理解されたい。これらの特定の特徴および行為は、本発明を実現する例示的な形態として開示したものである。

【図面の簡単な説明】

【図1】サーバレス分散ファイルシステムをサポートする例示的なネットワーク環境を示

10

20

30

40

50

す図である。

【図2】分散ファイルシステムに参加している図1の装置のうちの1台を表す例示的なコンピューティング装置の論理的構成要素を示す図である。

【図3】図1の分散ファイルシステムを実現するのに使用される、より一般的なコンピュータ環境を示す図である。

【図4】複数のサブツリーを有する名前空間ルートを含む、例示的な階層的名前空間を示す図である。

【図5】サブツリーの管理責任を他のディレクトリグループに委任する例示的なプロセスを示す流れ図である。

【図6】特定のパス名の管理に責任を負うディレクトリグループを探索する例示的なプロセスを示す流れ図である。

10

【図7】サーバレス分散ファイルシステム内でのファイルおよび対応するディレクトリエントリの例示的な記憶を示す図である。

【図8】サーバレス分散ファイルシステム内にファイルを記憶する例示的なプロセスを示す流れ図である。

【図9】特定のオブジェクトのオープンを許可するか否かを決定する例示的なプロセスを示す流れ図である。

【符号の説明】

100 ネットワーク環境

102, 104, 106, 108 クライアントコンピューティング装置

20

110 データ通信ネットワーク

120, 124, 128 分散システム部分

122, 126, 130, 132 ローカル部分

150 分散ファイルシステム

700 ファイルシステム

702, 704, 706, 708, 710, 712, 714, 716, 718, 720,

722, 724 コンピュータ

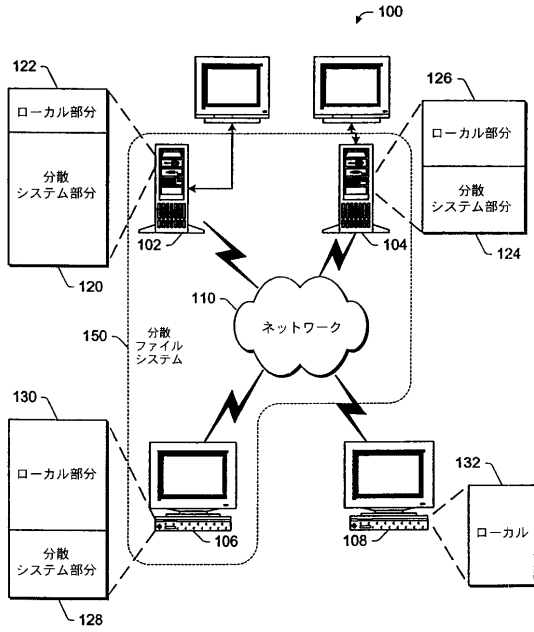
726 ビザンチングループ

728 ファイル

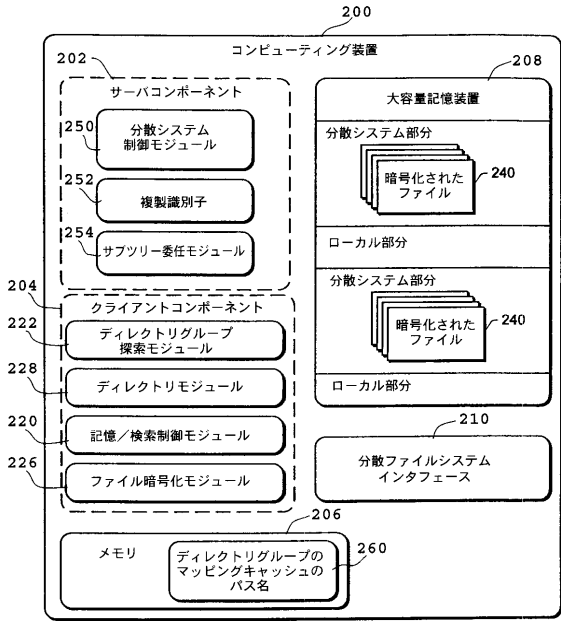
730 ディレクトリエントリ

30

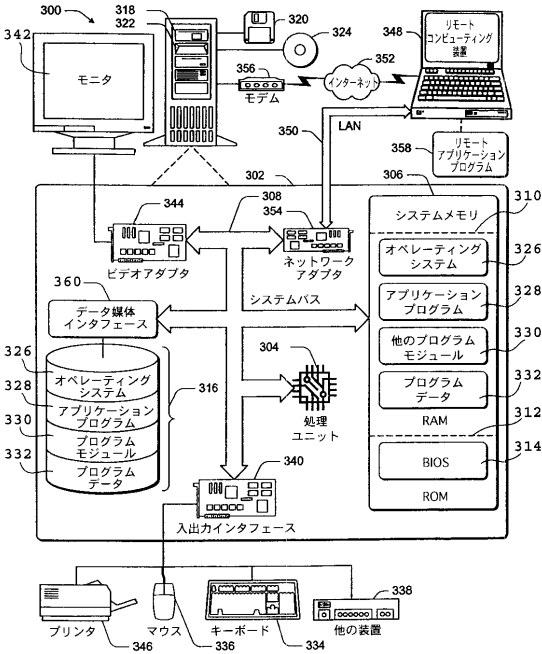
【図1】



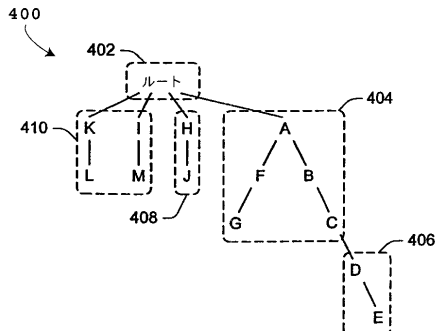
【図2】



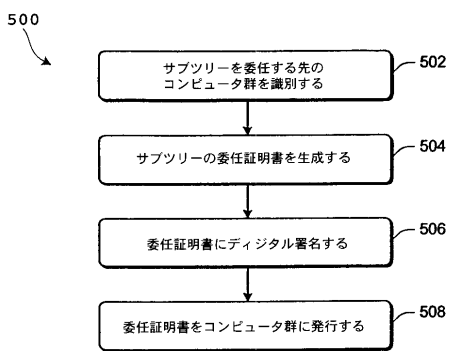
【図3】



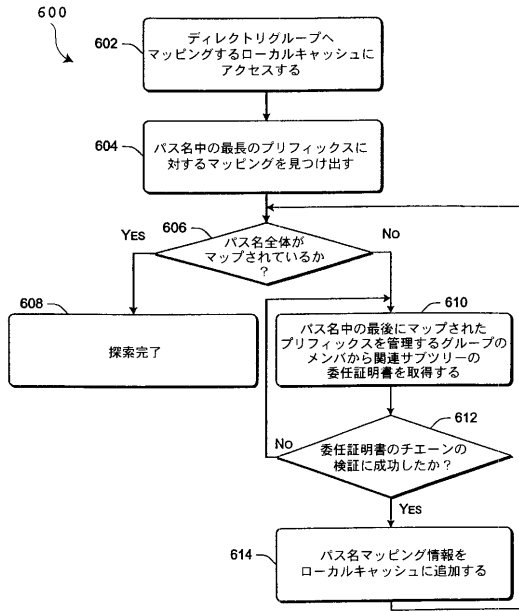
【図4】



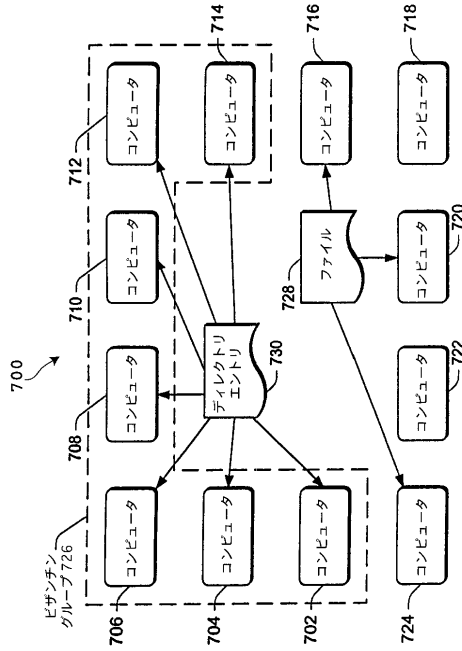
【図5】



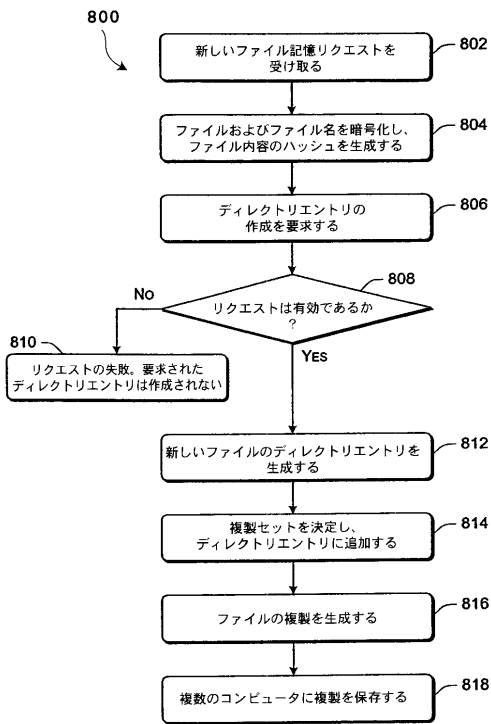
【図6】



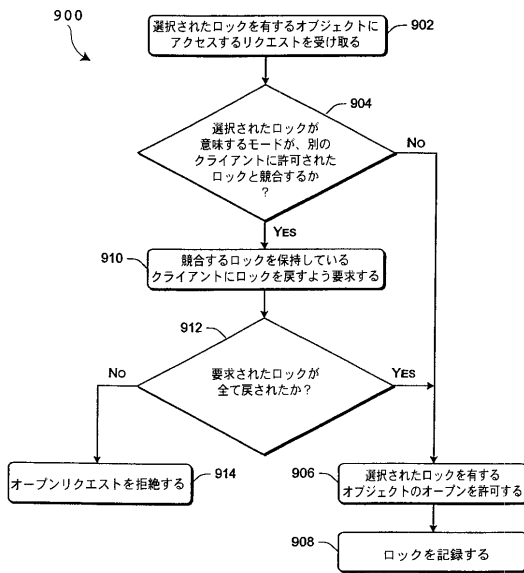
【図7】



【図8】



【図9】



フロントページの続き

- (72)発明者 ウィリアム ジェイ . ボロスキー
アメリカ合衆国 98027 ワシントン州 イサコア サウスイースト ミラーモント ドライ
ブ 24622
- (72)発明者 ジェラルド サーマック
アメリカ合衆国 98011 ワシントン州 ボズエル ノースイースト 141 ストリート
8908
- (72)発明者 ジョン アール . デューサー
アメリカ合衆国 98007 ワシントン州 ベルビュー ノースイースト 16 ストリート
14705
- (72)発明者 マービン エム . セイマー
アメリカ合衆国 98006 ワシントン州 ベルビュー 137 アベニュー サウスイースト
4440
- (72)発明者 ロバート ピー . ワッテンホファー
スイス シーエイチ - 8006 チューリッヒ クラウジウスシュトラーセ 35

審査官 工藤 嘉晃

- (56)参考文献 GEORGE COULOURIS 他 著 / 水野 忠則 他 訳 , 分散システム , 日本 , 株式会社電気書院 , 1997年10月31日 , 第2版 , p.785-p.827

(58)調査した分野(Int.Cl. , D B名)

G06F 12/00

JSTPlus(JDreamII)