

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2005-539315

(P2005-539315A)

(43) 公表日 平成17年12月22日(2005.12.22)

(51) Int.Cl. ⁷	F I	テーマコード (参考)
G06F 12/00	G06F 12/00	546K
G06F 13/00	G06F 13/00	540B
		5B082

審査請求 未請求 予備審査請求 未請求 (全 39 頁)

(21) 出願番号 特願2004-536622 (P2004-536622)
 (86) (22) 出願日 平成15年9月16日 (2003.9.16)
 (85) 翻訳文提出日 平成17年4月28日 (2005.4.28)
 (86) 国際出願番号 PCT/US2003/029398
 (87) 国際公開番号 W02004/025429
 (87) 国際公開日 平成16年3月25日 (2004.3.25)
 (31) 優先権主張番号 10/245,828
 (32) 優先日 平成14年9月16日 (2002.9.16)
 (33) 優先権主張国 米国 (US)

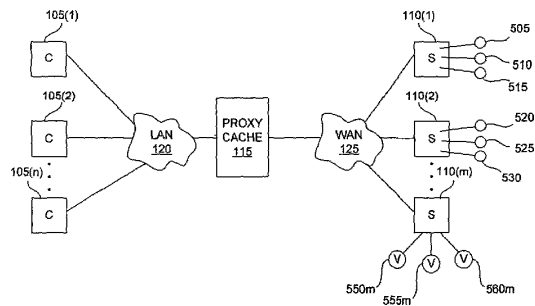
(71) 出願人 500261341
 ネットワーク・アプライアンス・インコー
 ポレイテッド
 アメリカ合衆国94089カリフォルニア
 州サニーベイル、イースト・ジャーバ・ド
 ライブ495番
 (74) 代理人 100087642
 弁理士 古谷 聡
 (74) 代理人 100076680
 弁理士 溝部 孝彦
 (74) 代理人 100121061
 弁理士 西山 清春

最終頁に続く

(54) 【発明の名称】 プロキシ・キャッシュに関する装置および方法

(57) 【要約】

一実施形態において本発明は、データをネットワーク内にキャッシュする装置を提供する。この装置は、クライアント(105)からオブジェクトの要求を受信し、サーバ(110)からデータブロックをフェッチするように構成されたプロキシ・キャッシュ(115)を含む。プロキシ・キャッシュは、データブロックをオブジェクト内に階層関係でキャッシュするように構成される。オブジェクトには、例えばデータファイルやディレクトリがある。プロキシ・キャッシュにキャッシュされたデータブロックは、クライアントからの要求に基づくアクティブデータセットを規定する。



【特許請求の範囲】**【請求項 1】**

データをネットワーク内にキャッシュする装置であって、クライアントからオブジェクトの要求を受信し、サーバからデータブロックをフェッチするように構成されたプロキシ・キャッシュを含み、

前記プロキシ・キャッシュは、階層関係の前記データブロックをアクティブデータセットとしてオブジェクト内にキャッシュするように構成される、装置。

【請求項 2】

前記データブロックがバッファ・ツリーに入れられ、ファイル内に前記階層関係が可能である、請求項 1 に記載の装置。

10

【請求項 3】

クライアントからの要求に応答して、データブロックが前記アクティブデータセットに配置される、請求項 1 に記載の装置。

【請求項 4】

前記アクティブデータセットは、少なくともサーバに格納されたデータの一部を含む、請求項 1 に記載の装置。

【請求項 5】

前記アクティブデータセットは、特定の時間によって定義される、請求項 1 に記載の装置。

【請求項 6】

前記プロキシ・キャッシュは、サーバによって要求されたファイルがデータブロックを欠いていた場合、クライアントからの要求に応答してサーバから前記データブロックをフェッチするように構成される、請求項 1 に記載の装置。

20

【請求項 7】

前記プロキシ・キャッシュは、サーバによって要求されたファイルがプロキシ・キャッシュ内に存在しなかった場合、クライアントからの要求に応答してサーバから前記データブロックをフェッチするように構成される、請求項 1 に記載の装置。

【請求項 8】

前記プロキシ・キャッシュと前記クライアントが、オープン・スタンダード・プロトコルを用いて通信する、請求項 1 に記載の装置。

30

【請求項 9】

前記プロトコルは、インターネット・プロトコルを含む、請求項 8 に記載の装置。

【請求項 10】

前記プロトコルは、ネットワーク・ファイル・システム (NFS) プロトコルを含む、請求項 8 に記載の装置。

【請求項 11】

前記プロトコルは、コモン・インターネット・ファイル・システム (CIFS) プロトコルを含む、請求項 8 に記載の装置。

【請求項 12】

前記プロキシ・キャッシュは、クライアントによって要求されたデータファイルにおけるデータブロックの有無を判定するための計算の一部として重複しない名前を計算する、請求項 1 に記載の装置。

40

【請求項 13】

置換ポリシーに基づいて、前記プロキシ・キャッシュから前記アクティブデータセットの一部を消去することができる、請求項 1 に記載の装置。

【請求項 14】

前記オブジェクトはデータファイルを含む、請求項 1 に記載の装置。

【請求項 15】

オブジェクトはファイルの一部を含む、請求項 1 に記載の装置。

【請求項 16】

50

データをネットワーク内にキャッシュする方法であって、

クライアントからオブジェクトの要求を受信するステップと、

要求された前記オブジェクト中の各データブロックがキャッシュ内にあれば、要求されたオブジェクトをクライアントに送信するステップと、

要求された前記オブジェクト中のデータブロックが欠けていた場合、該データブロックをサーバからフェッチし、該データブロックをアクティブデータセットの一部として階層関係に埋め込み、要求された前記オブジェクトをクライアントに送信するステップと、
からなる方法。

【請求項 17】

サーバによって要求されたファイルが存在しなかった場合に、クライアントからサーバへの要求に
10 応答して、サーバからデータブロックをフェッチするステップをさらに含む、請求項 16 に記載の方法。

【請求項 18】

前記データブロックがバッファ・ツリーに入れられ、ファイル内に階層関係が可能である、請求項 16 に記載の方法。

【請求項 19】

前記アクティブデータセットは、少なくともサーバに格納されたデータの一部を含む、請求項 16 に記載の方法。

【請求項 20】

前記アクティブデータセットは、特定の時間によって定義される、請求項 16 に記載の方法。
20

【請求項 21】

前記オブジェクトの要求は、オープン・スタンダード・プロトコルを用いて送信される、請求項 16 に記載の方法。

【請求項 22】

前記プロトコルは、インターネット・プロトコルを含む、請求項 21 に記載の方法。

【請求項 23】

前記プロトコルは、ネットワーク・ファイル・システム (NFS) プロトコルを含む、請求項 21 に記載の方法。

【請求項 24】

前記プロトコルは、コモン・インターネット・ファイル・システム (CIFS) プロトコルを含む、請求項 21 に記載の方法。
30

【請求項 25】

クライアントによって要求されたデータファイル中のデータブロックの有無を判定する計算の一部として重複しない名前を計算するステップをさらに含む、請求項 16 に記載の方法。

【請求項 26】

置換ポリシーに基づいて前記アクティブデータセットの一部を消去するステップをさらに含む、請求項 16 に記載の方法。

【請求項 27】

クライアントからオブジェクトの要求を受信する命令と、
要求された前記オブジェクトの各データブロックがキャッシュ内に存在する場合、要求された前記オブジェクトを前記クライアントに送信する命令と、
要求された前記オブジェクト中のデータブロックが欠けていた場合、該データブロックをサーバからフェッチし、該データブロックをアクティブデータセットとして階層関係に埋め込み、要求された前記オブジェクトをクライアントに送信する命令と、
40 が記録された機械読取可能媒体を含む製品。

【請求項 28】

データをネットワーク内にキャッシュする装置であって、

クライアントからオブジェクトの要求を受信するための受信手段と、
50

前記受信手段に接続され、要求された前記オブジェクトの各データブロックがキャッシュ内に存在する場合、要求された前記オブジェクトを前記クライアントに送信する手段と、

前記受信手段に接続され、要求された前記オブジェクト中のデータブロックが欠けていた場合、該データブロックをサーバからフェッチし、該データブロックをアクティブデータセットの一部として階層関係に埋め込み、要求された前記オブジェクトを前記クライアントに送信する手段と、

からなる装置。

【請求項 29】

データをネットワーク内にキャッシュする装置であって、

クライアントからオブジェクトの要求を受信し、サーバからデータブロックをフェッチするように構成されたプロキシ・キャッシュを含み、

前記プロキシ・キャッシュは、階層関係の前記前記データブロックをアクティブデータセットとしてオブジェクト内にキャッシュするように構成され、該アクティブデータセットが、前記クライアントからの要求に基づいて定義される、装置。

【請求項 30】

前記アクティブデータセットはデータ置換法を用いて変更される、請求項 29 に記載の装置。

【請求項 31】

前記データ置換法は、先入れ先出し法、最長時間未使用法、データ重要度法、および、ユーザ入力を利用したデータ置換法のうちのいずれか 1 つである、請求項 29 に記載の装置。

【請求項 32】

前記プロキシ・キャッシュに接続された全てのクライアントが前記アクティブデータセットを相互共有できるようにするために、前記アクティブデータセットは、前記プロキシ・キャッシュに接続された少なくとももう 1 つのクライアントから利用することができる、請求項 29 に記載の装置。

【請求項 33】

前記プロキシ・キャッシュは、ファイルハンドルを変換するように構成され、該ファイルハンドルが特定のサーバにマッピングされていることを示すようにする、請求項 29 に記載の装置。

【請求項 34】

前記ファイルハンドルは、ファイルシステム識別子およびサーバ・インターネット・プロトコル・アドレスを利用してハッシュアルゴリズムを用いて変換される、請求項 29 に記載の装置。

【請求項 35】

データをネットワーク内にキャッシュする方法であって、

クライアントからオブジェクトの要求を受信するステップと、

要求された前記オブジェクト中の少なくとも 1 つのデータブロックがキャッシュ内に存在するか否かを判定するステップと、

要求された前記オブジェクト中のデータブロックが欠けていた場合、該データブロックをサーバからフェッチし、該データブロックをキャッシュに割り当て、要求された前記オブジェクトを前記クライアントに送信するステップと、

からなる方法。

【請求項 36】

各データブロックが、オブジェクト内の階層関係に入れられる、請求項 35 に記載の方法。

【請求項 37】

前記アクティブデータセットは、少なくとも前記サーバに格納されたデータの一部を含む、請求項 35 に記載の方法。

10

20

30

40

50

- 【請求項 38】
前記アクティブデータセットは、特定の時間によって定義される、請求項 35 に記載の方法。
- 【請求項 39】
前記オブジェクトの要求は、オープン・スタンダード・プロトコルを用いて送信される、請求項 35 に記載の方法。
- 【請求項 40】
前記プロトコルは、インターネット・プロトコルを含む、請求項 39 に記載の方法。
- 【請求項 41】
前記プロトコルは、ネットワーク・ファイル・システム (NFS) プロトコルを含む、請求項 40 に記載の方法。 10
- 【請求項 42】
前記プロトコルは、コモン・インターネット・ファイル・システム (CIFS) プロトコルを含む、請求項 40 に記載の方法。
- 【請求項 43】
クライアントによって要求されたデータファイル中のデータブロックの有無を判定する計算の一部として重複しない名前を計算するステップをさらに含む、請求項 35 に記載の方法。
- 【請求項 44】
置換ポリシーに基づいて前記アクティブデータセットの一部を消去するステップをさらに含む、請求項 35 に記載の方法。 20
- 【請求項 45】
前記データ置換法は、先入れ先出し法、最長時間未使用法、データ重要度法、および、ユーザ入力を利用したデータ置換法のうちのいずれか 1 つである、請求項 44 に記載の方法。
- 【請求項 46】
前記プロキシ・キャッシュに接続された全てのクライアントが前記アクティブデータセットを相互共有できるようにするために、前記アクティブデータセットは、前記プロキシ・キャッシュに接続された少なくとももう 1 つのクライアントから利用することができる、請求項 35 に記載の方法。 30
- 【請求項 47】
前記プロキシ・キャッシュは、ファイルハンドルを変換するように構成され、該ファイルハンドルが特定のサーバにマッピングされていることを示すようにする、請求項 35 に記載の方法。
- 【請求項 48】
前記ファイルハンドルは、ファイルシステム識別子およびサーバ・インターネット・プロトコル・アドレスを利用してハッシュアルゴリズムを用いて変換される、請求項 47 に記載の方法。
- 【請求項 49】
前記プロキシ・キャッシュは、クライアントからの要求に回答して重複しないファイル名を計算し、該重複しないファイル名に基づいて、要求された前記オブジェクトの各データブロックが前記プロキシ・キャッシュに格納されているか否かを判定する一致キャッシュファイル識別子を得る、請求項 1 に記載の装置。 40
- 【請求項 50】
前記一致ローカルキャッシュファイル識別子は、適当な inode レコードを指し示す inode ファイルへの索引となる、請求項 49 に記載の装置。
- 【請求項 51】
前記 inode レコードは、要求された前記オブジェクトのデータブロックが前記プロキシ・キャッシュに格納されているか否かを示す inode 番号を含む、請求項 50 に記載の装置。 50

【請求項 5 2】

要求された前記オブジェクト中の関連データブロックが前記プロキシ・キャッシュ内に格納されていない場合、`inode` 番号は特別な値 X をとり、該特別な値 X は、 -1 、 -3 、または、その他の値などの値のうちの一つである、請求項 5 1 に記載の装置。

【請求項 5 3】

要求された前記オブジェクト中の関連データブロックが前記プロキシ・キャッシュ内に格納されている場合、`inode` 番号は特別な値 X にならない、請求項 5 1 に記載の装置。

【請求項 5 4】

前記 `inode` ファイルは、オブジェクトに関するサーバ側属性を含む第 2 の `inode` レコードをさらに指し示す、請求項 4 9 に記載の装置。 10

【請求項 5 5】

前記重複しないファイル名に基づく一致キャッシュファイル識別子が前記プロキシ・キャッシュ内にない場合、要求された前記オブジェクトのデータブロックは、前記プロキシ・キャッシュ内に格納されていない、請求項 4 9 に記載の装置。

【請求項 5 6】

前記クライアントからの要求に応答して重複しないファイル名を計算するステップと、前記重複しないファイル名に基づいて、要求された前記オブジェクトの各データブロックがプロキシ・キャッシュに格納されているか否かを判定する一致キャッシュファイル識別子を得るステップと、
をさらに含む、請求項 1 6 に記載の方法。 20

【請求項 5 7】

前記一致ローカルキャッシュファイル識別子は、適当な `inode` レコードを指し示す `inode` ファイルへの索引となる、請求項 5 6 に記載の方法。

【請求項 5 8】

前記 `inode` レコードは、要求された前記オブジェクトのデータブロックが前記プロキシ・キャッシュに格納されているか否かを示す `inode` 番号を含む、請求項 5 7 に記載の方法。

【請求項 5 9】

要求された前記オブジェクト中の関連データブロックが前記プロキシ・キャッシュ内に格納されていない場合、`inode` 番号は特別な値 X をとり、該特別な値 X は、 -1 、 -3 、または、その他の値などの値のうちの一つである、請求項 5 8 に記載の方法。 30

【請求項 6 0】

要求された前記オブジェクト中の関連データブロックが前記プロキシ・キャッシュ内に格納されている場合、`inode` 番号は特別な値 X にならない、請求項 5 8 に記載の方法。

【請求項 6 1】

前記 `inode` ファイルは、オブジェクトに関するサーバ側属性を含む第 2 の `inode` レコードをさらに指し示す、請求項 5 6 に記載の方法。

【請求項 6 2】

前記重複しないファイル名に基づく一致キャッシュファイル識別子が前記プロキシ・キャッシュ内にない場合、要求された前記オブジェクトのデータブロックは、前記プロキシ・キャッシュ内に格納されていない、請求項 5 6 に記載の方法。 40

【請求項 6 3】

前記オブジェクトはディレクトリを含む、請求項 1 に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は概してコンピュータ・ネットワークに関する。詳しくは、本発明の実施形態は概してファイルシステムおよび記憶装置に関する。

【背景技術】

【0002】

大きな組織は地理的に分散配置されている。それらは、高速のローカル・エリア・ネットワーク（LAN）を備えた少数の大型中央サイトを有し、ローカルデータを中央エリアまたは中央システムにまとめて整理して置くことができる。しかしながら、大きな組織は、小さな遠方の事務所、支社、および/または、他の末端位置をさらに有する場合があります。それらから中央集中化されたデータ保管場所までの接続の帯域幅は比較的狭く、遅延も比較的大きい。末端位置には、例えばサーバのバックアップや故障したハードウェアおよび/またはソフトウェアの交換といったローカルデータのニーズに携わるIT（情報技術）スタッフが居ない場合がある。遠方の事務所に居るユーザは、中央サイトに居る相手方に比べて、企業ITインフラストラクチャの能力、信頼性、および、サービスが非常に低いことを知っている。

10

【0003】

大型サイト間におけるデータの共有は厄介でもある。大型サイトは通常、専用の広帯域リンクによって相互接続される。しかしながら、サイト間の地理的距離は、多くの用途にとって許容されない遅延を発生させる。

【0004】

現在の解決策は、各末端位置にファイルサーバ（ネットワーク・アプライアンス・インコーポレイテッドから市販されている「ファイラ」など）を設置し、サーバをワイド・エリア・ネットワーク（WAN）を介して中央位置に自動的にバックアップすることである。サイト間で共有される読出し専用データは、同期ミラーリングによって複製される。ハイパーテキスト・トランスファ・プロトコル（HTTP）およびストリーミング・トラフィックが、例えばネットワーク・アプライアンス・インコーポレイテッドから市販されている「NetCache」のようなネットキャッシュを用いてキャッシュされる。

20

【0005】

末端位置におけるファイルサーバの管理は、コストおよび/またはリソースの点で高くつく可能性がある。WANを介したデータのバックアップ作業は、慎重な計画と管理を必要とする。WANを介して読出し専用ボリュームを複製することは、やり過ぎであることが多い。多くの場合、実際に日常的に使用されるのは、ボリュームのうちの小さな割合の部分（/usr/local/binなど）だけである。一般に、動的データセット（ユーザのホームディレクトリなど）を複製することはなく、動的データセットは各位置で個別に管理しなければならない。

30

【0006】

従って、上記の製品および方法は、特定の能力および機能に制限され、数々の制約を受ける。

【発明の開示】

【課題を解決するための手段】

【0007】

本発明は一実施形態において、データをネットワーク内にキャッシュする装置を提供する。この装置は、クライアントからオブジェクトの要求を受信し、サーバからデータブロックをフェッチするように構成されたプロキシ・キャッシュを含む。プロキシ・キャッシュは、データブロックをオブジェクト内に階層関係でキャッシュするように構成される。オブジェクトには、例えばデータファイル、ディレクトリ、データブロックのようなデータファイルの一部などがある。プロキシ・キャッシュにキャッシュされたデータブロックは、クライアントからの要求に基づいてアクティブデータセットを定義する。

40

【0008】

それらの対策、並びに、説明の進行につれて当業者には明らかとなるであろう付随的な対策および特徴は、添付の図面に例として符号付きで示される装置、アセンブリ、本発明

50

の実施形態の方法、およびそれらの好ましい実施形態によって達成される。

【発明を実施するための最良の形態】

【0009】

図面を参照し、本発明の種々の実施形態について説明する。特に断わりがない限り、全部図面を通じて、同じ符号は同じ部品を意味している。

【0010】

本発明の実施形態を完全に理解してもらうために、本明細書には、構成要素および/または方法の例など、多数の具体的詳細を記載している。しかしながら、本発明が、それらの具体的詳細のうち1以上が無くても実施できるものであり、他の装置、システム、方法、構成要素、材料、部品などを用いて実施することも可能であることは、当業者には明らからであろう。本発明の実施形態が不明確になることを避けるために、他の例では、周知の構造、材料、処理については図示も説明も詳しくしていない。

10

【0011】

図1Aは、クライアント装置105(1)~105(n)、サーバ110(1)~110(m)およびプロキシ・キャッシュ(すなわちプロキシ・アプライアンス)115を含むネットワーク100を示すブロック図である。これらの構成要素は、本発明の一実施形態に従ってフォワード・プロキシ構成を形成している。以下で説明するように、本発明の他の実施形態は、リバース・プロキシ構成の少なくとも1つのプロキシ・キャッシュを含むネットワークを有する。本発明の実施形態は、下記の機能のうち少なくとも一部を有する。すなわち、提供する。

20

- (1) 位置に非依存、および、名前空間の統合
- (2) オン・デマンドでまばらな整合のとれたデータの複製
- (3) ロード・バランシング
- (4) 遠隔切断アクセスおよびデータ変更
- (5) プロトコル変換

これらの仮想化を適用することにより、従来の方法に関連する法外なコストを伴うことなく、分散型ストレージの基盤を構築することが可能になる。

【0012】

上記の変数nおよびmは任意の適当な整数値である。従って、クライアント装置(本明細書では通常、クライアント105と呼ぶ)の数と、サーバ(本明細書では通常、サーバ110と呼ぶ)の数は、変更してもよい。例えば、ネットワーク100は、たった1つのクライアント装置105、および/または、たった1つのサーバ110を用いて実施することもできる。クライアント装置105は、ローカル・エリア・ネットワーク(LAN)120を介してプロキシ・キャッシュ115に接続され、サーバ110は、ワイド・エリア・ネットワーク(WAN)125を介してプロキシ・キャッシュ115に接続される。

30

【0013】

図1Aに示すフォワード・プロキシ構成によれば、データをプロキシ・キャッシュ115にオン・デマンドでキャッシュすることができる。つまり、プロキシ・キャッシュ115によるオン・デマンドでまばらな整合のとれたデータの複製が可能となる。このオン・デマンド・キャッシュ処理は、従来複製方法に比べて効率が高く、従来複製方法で使用されているような特殊なソフトウェアを必要としないという利点がある。

40

【0014】

一実施形態において、プロキシ・キャッシュ115とクライアント装置105は、ネットワーク・ファイル・システム(NFS)プロトコルのようなオープン・スタンダード・プロトコル130を用い、LAN120を介して互いに通信することができる。以下で説明するように、オープン・スタンダード・プロトコル130は、コモン・インターネット・ファイル・システム(CIFS)プロトコルのような他の適当なオープン・スタンダード・プロトコルであってもよい。プロキシ・キャッシュ115とサーバ110は、NFSのようなオープン・スタンダード・プロトコル135を用い、ワイド・エリア・ネットワーク(WAN)125を介して互いに通信することができる。本発明の一実施形態は、W

50

A N 1 2 5 および L A N 1 2 0 にオープン・スタンダード・プロトコルを使用しているの
で、従来の方法で必要とされる特殊なソフトウェアが不要であるという利点を有する。特
に、N F S のようなオープン・スタンダード・プロトコルを使用することにより、クライ
アント 1 0 5 として異種のクライアントを使用することが可能となる。言い換えれば、「
異種のクライアント」という用語は、異なる製造業者またはベンダーで製造された種々の
クライアントを都合よく使用出来ることを意味している。クライアントは、サーバ 1 1 0
と通信するための特殊なソフトウェアを必要としない。

【 0 0 1 5 】

さらに、オープン・スタンダード・プロトコルを使用すると、プロキシ・キャッシュ 1
1 5 をネットワーク 1 0 0 内に設置するとき通常ならば必要となるクライアント装置 1
0 5 の設定変更も不要になる。従って、プロキシ・キャッシュ 1 1 5 を使用すると、ネッ
トワーク管理に関するオーバヘッドおよびコストを低減することができる。なお、プロキ
シ・キャッシュ 1 1 5 は、ネットワーク 1 0 0 内のサーバ 1 1 0 やクライアント 1 0 5 か
ら離れた場所から管理したり、それらから離れた場所に設置してもよい。

10

【 0 0 1 6 】

通常、プロキシ・キャッシュ 1 1 5 は、プロキシ・キャッシュ 1 1 5 に接続されたサー
バ 1 1 0 およびクライアント 1 0 5 を識別するとともに、サーバ 1 1 0 に格納された特定
のデータファイルを識別することができる。

【 0 0 1 7 】

例えば、クライアント 1 0 5 がデータを要求したときに、そのデータがプロキシ・キャ
ッシュ 1 1 5 になければ、「キャッシュ・ミス」が発生する。一方、要求されたデータが
プロキシ・キャッシュ 1 1 5 にあれば、「キャッシュ・ヒット」が発生する。これらの処
理については、図 1 B および図 1 C を参照して後で詳しく説明する。

20

【 0 0 1 8 】

キャッシュ・ヒットが発生すると、プロキシ・キャッシュ 1 1 5 は、要求されたデータ
を要求元クライアント 1 0 5 に送信する。キャッシュ・ミスが発生すると、プロキシ・キャ
ッシュ 1 1 5 は、データをサーバ 1 1 0 に要求し、次いでそのデータを要求元クライ
アント 1 0 5 に渡す。サーバ 1 1 0 から得られたデータは、アクティブデータセット 1 4 0
としてプロキシ・キャッシュ 1 1 5 内にキャッシュされ、他のクライアント 1 0 5 から直
ぐに利用できるようにされる。アクティブデータセットは、クライアントおよびクライ
アントのアプリケーションの属性であり、所与の時間（すなわち時間窓）内に参照され、キャ
ッシュされたデータである。従って、アクティブデータセットは、所与の時間に
応じて異なる場合がある（例えば、所与の時間を 2 分としたときのアクティブデータ
セットは、所与の時間を 1 日または 1 週間としたときのアクティブデータセットとは異なる場合がある）。

30

【 0 0 1 9 】

プロキシ・キャッシュ 1 1 5 は、キャッシュした複製データを複数のクライアント 1 0
5 間で相互共有させることができる。この相互共有状態は、クライアント 1 0 5 からサー
バ 1 1 0 に要求された特定のデータが、少なくとも 1 つの他のクライアント 1 0 5 によ
ってさらに要求される可能性が高いことを前提としている。

40

【 0 0 2 0 】

一実施形態において、プロキシ・キャッシュ 1 1 5 はアクティブデータセット 1 4 0 を
キャッシュする。アクティブデータセット 1 4 0 とは、最近要求されたデータセットまた
は要求頻度の高いデータセットであって、プロキシ・キャッシュ 1 1 5 からまだ消去され
ていないデータセットである。アクティブデータセット 1 4 0 は、サーバ（複数の場合も
あり）1 1 0 に格納されているデータのサブセットである。クライアント 1 0 5 がサーバ
1 1 0 からデータファイルの特定のコピー 1 5 0 を読み出すために読み出し要求 1 4 5 を送
ると、その読み出し要求 1 4 5 はプロキシ・キャッシュ 1 1 5 に受信され、そこでそのファ
イルまたはフォルダの特定部分（例えば、データブロックなど）が、プロキシ・キャッシ
ュ 1 1 5 内にローカルにキャッシュされているか否かがチェックされる。要求された特定

50

のデータファイルがアクティブデータセット140内にあった場合、キャッシュ・ヒット状態が発生し、プロキシ・キャッシュ115は、要求されたデータをLAN120を介して要求元クライアント105に送信する。

【0021】

一方、要求されたデータファイル中の或るデータブロックがアクティブデータセットの一部として格納されていなかった場合、キャッシュ・ミス状態が発生する。すると、プロキシ・キャッシュ115は、要求155をサーバ110に送り、要求されたデータファイル中の欠けているデータブロックのコピーを得る。

【0022】

一実施形態において、キャッシュされたファイルは、アクティブデータセット140内でそのファイルが適当な置換方法（例えば、データの先入れ先出し管理や、最長時間未使用（LRU）アルゴリズムなど）によって置換されるまで、アクティブデータセット140の一部として保持される。当然ながら、以下で説明するように、一実施形態において、アクティブデータセット140の管理に関しては、他のデータ置換方法を使用してもよい。例えば、あるファイルを永久記憶用として指定し、クライアント105のユーザがアクティブデータセット140からそのファイルを抹消（消去）するコマンドを発行するまで、そのファイルを保持するようにしてもよい。

【0023】

図1Bは、本発明の一実施形態によるプロキシ・キャッシュ115を示すブロック図である。プロキシ・キャッシュ115は、IPを利用したネットワーク・トラフィック（または、ファイバ・チャネルやストレージ・エリア・ネットワークなどに関する他のタイプのトラフィック）を構文解析するための、ネットワーク・インタフェース161、ネーミング仮想化層162、NFSサービス163、および/または、CIFSサービス164、ローカル・アクセス層165、NFSフィルエンジン166、および/または、CIFSフィルエンジン167、および/または、他の適当なタイプのフィルエンジン、リモート・フィル層168、ファイルシステム層169（例えば、Write-Anywhere-File-Layout、すなわちWAF Lなど）、ストレージデバイス・マネージャ170（Redundant Array Of Independent (or Inexpensive) Disks層、すなわちRAID層など）、および、ストレージ・ディスク（複数の場合もあり）171を含む。

【0024】

ネットワーク・インタフェース161は、クライアント105からストレージ関係のサービス要求を受信するためのコンポーネントを含む。

【0025】

一般にファイルシステムは、情報をディレクトリとファイルの階層構造としてストレージデバイス（例えば、ディスクなど）上に論理編成することができる。（ディスク上の）各ファイルは、テキストなどの情報を格納するように構成されたディスクブロックの集合として実施されるのに対し、ディレクトリは、他のファイルやディレクトリに関する情報が格納された特殊形式のファイルとして実施される。

【0026】

ストレージデバイス・マネージャ170は、ストレージシステム内のストレージデバイス171を管理する。ストレージデバイス・マネージャ170は、ファイルシステム169から読み出しコマンドおよび書き込みコマンドを受信し、それらのコマンドを処理することによって、ストレージシステムにアクセスする。ストレージデバイス・マネージャ170は、ファイルシステム169からブロックの論理アドレスを得て、その論理アドレスをストレージシステム内の1以上のストレージデバイス171における物理アドレスに変換する。一実施形態において、ストレージデバイス・マネージャ170は、ストレージデバイスをRAID（Redundant Array of Independent, or Inexpensive, Disks）に従って管理する。

【0027】

一般に、ディスク・ストレージは大抵、記憶空間の全体的論理構成を定義している物理

10

20

30

40

50

ストレージディスクによって形成された1以上のストレージ「ボリューム」として実施される。各ボリュームは通常、そのボリューム独自のファイルシステムに関連しているので、「ボリューム」という用語と「ファイルシステム」という用語は、同じ意味で使用される。ボリューム内のディスクは通常、1以上のRAIDグループとして編成される。

【0028】

図1Bに示す他のモジュールの機能については、図1Cを参照して説明する。

【0029】

キャッシュ・ヒット

次に、図1Bおよび図1Cのブロック図を参照し、本発明の一実施形態によるプロキシ・キャッシュ115の動作について説明する。ファイルハンドルは、ファイルシステム内のファイル等のオブジェクトに名前を付けるのに使用される。ファイルハンドルについては、図4を参照して後で詳しく説明する。

【0030】

プロキシ・キャッシュ115(図1B)が、クライアント105から、例えばFH = 「FILE1 ON CACHE」の値を有するファイルハンドル180を含む読出し要求を受信したものと仮定する。すると、ネットワーク・インタフェース161はその要求をネーミング仮想化層162に転送し、仮想化層が「FILE1」を、例えばIP(インターネット・プロトコル)アドレス10.56.20.34を有するサーバ110上のサーバ・ファイルハンドルFH = 「FILE1 ON SERVER1」にマッピングする。ローカル・アクセス層165は、ファイルサーバID(FSid)値181およびサーバIPアドレス10.56.20.34に基づいて重複しない(一意の)名前182を計算する。一実施形態において、重複しない名前を計算するアルゴリズムはMD5ハッシュアルゴリズムである。MD5ハッシュアルゴリズムはメッセージを受け取ると、それをメッセージ・ダイジェストと呼ばれる固定桁数の文字列に変換する。また、要求172がNFS要求であるかCIFS要求であるかに応じて、NFSサービス層163またはCIFSサービス層164が、その要求172を構文解析する機能を提供する。

【0031】

ファイルシステム層169は、重複しない名前182に基づいて検索機能を実施し、ローカル・キャッシュファイルID185を得る。一実施形態において、この検索機能は、ファイルシステム内のモジュール184によって実施される。一致するローカル・キャッシュファイルIDが無ければ、キャッシュ・ミスが発生する。キャッシュ・ミスについては、後で詳しく説明する。

【0032】

ローカル・キャッシュファイルID183は、適当なinodeレコード186を指し示す、inodeファイル185への索引である。inodeファイル105のinodeレコードは、所与のファイルシステムに関連するinodeファイルを表わす情報を有する。一般に、inodeレコードは、ファイルに関するメタデータ(属性)などの情報の格納に使用されるデータ構造であり、ファイル・データブロックは、そのファイルの実際のデータの格納に使用される構造である。inodeレコードに格納される情報には、例えば、ファイルの所有者、ファイルのアクセス・パーミッション、ファイルのサイズ、ファイルタイプ、ファイルのデータブロックのディスク上の位置に対する基準などがある。inodeファイル185のinodeレコード186は、まとめて符号188で示すような、ファイルシステム・データブロック(例えば、WAFLEデータブロックなど)へのポインタを有している。ファイル・データブロックは、ファイルシステムによって処理されるデータの最小アドレス指定可能量として定義される。ファイル・データブロックは、例えば4キロバイト(KB)のデータを格納する容量を有する。inodeレコード186は間接ブロックを指し示すことがあり、間接ブロックは他のファイル・データブロックまたは他の間接ブロックをさらに指し示すことがある。例えば、間接ブロック188dはブロック189a、189b、189cを指し示す場合があり、それらのブロックはそれぞれ、ファイル・データブロックまたは間接ブロックである。

10

20

30

40

50

【0033】

具体的には、`inode`レコード186内の`inode`番号187は、ファイル・データブロックまたは間接ブロックを指し示す。また、`inode`レコード186は、プロキシ・キャッシュ115にローカルに格納されたファイルの属性190をさらに含む。さらに、一実施形態において、ローカル・キャッシュファイルID183は、ファイルのサーバ側属性を含む適当な`inode`レコード193を指し示す、第2の`inode`ファイル192への索引でもある。サーバ側属性には、例えば、ファイルシステムID、ファイルID、ブロックサイズ、ハードリンクの数、ファイルシステム上の空き空間などがある。当然ながら、第2の`inode`ファイル192は、`inode`ファイル185に結合させてもよい。第2の`inode`ファイル192を作成する利点は、図1Cを参照して説明した機能の一部を実施するための第1の`inode`ファイル185の変更が、不要になることである。

10

【0034】

また、`inode`レコードは、ファイルシステム内のディスク上のファイルの表現に使用されるデータ構造であるファイルバッファ・ツリーを含む。図1Cでは、ブロック番号187およびブロック188, 189により、ファイルバッファ・ツリー191が形成されている。

【0035】

次に、IOベクトル194を形成する。IOベクトルは、ファイルバッファ・ツリー191のブロック番号187のリストである。ブロック番号187は、特定ファイルについてファイル・データブロックの有無を示す。一実施形態において、モジュール195はIOベクトル194を形成することができ、ファイル・データブロックの有無を判定する。一実施形態において、ブロック番号187は、ブロック番号187aに示すような特殊な値X（Xは例えば-1または-3をとりうる）を有する。この特殊な値Xは、プロキシ・キャッシュ115に格納されたローカルキャッシュファイルが、要求されたデータブロックを有していないことを示す値である。

20

【0036】

キャッシュ・ヒット状態になると、ブロック番号は特殊な値Xではなくなり、IOベクトルがストレージデバイス・マネージャ170（例えば、RAID層など）に送信される。ストレージデバイス・マネージャ170は、適当なストレージ・デバイス（複数の場合もあり）へのポインタを生成し、要求されたデータブロックをファイルシステムに送信するとともに、NFSサービス（NFS要求の場合）に送信する。そしてNFSサービスは、NFS応答173を要求元クライアント105に対して生成する。従って、要求元クライアント105は、要求したファイルを受信することができる。上記の方法は、ディレクトリ要求の処理に使用することもできる。

30

【0037】

ブロック番号リスト中に特殊な値Xを使用すると、要求されたファイルの欠けているデータブロックを有するプロキシ・キャッシュ（のストレージデバイス）におけるバッファの追跡が可能になり、従って散在するデータの追跡が可能になる。データブロックまたはキャッシュファイルはプロキシ・キャッシュ115から消去される場合があるので、（重複しない名前182を計算することによる）間接レベルによれば、特定の時刻に応じてファイルハンドルを`inode`ファイル内の異なるスロットを指し示すものに行うことができる。

40

【0038】

キャッシュ・ミス

IOベクトル194を作成した場合、ファイルバッファ・ツリー191のブロック番号が、要求されたファイル・データブロックが（ストレージデバイス171内の）バッファ内に無いことを意味する特殊な値Xであれば、キャッシュ・ミス状態が発生する。あるいは、重複しない名前182を計算した後、ファイルシステム169がテーブル検索機能を実施したときに一致するファイルID183が見付からなければ、キャッシュ・ミス状態

50

が発生する。

【0039】

NFS要求196のためのNFSフィルエンジン166（またはCIF S要求のためのCIF Sフィルエンジン）は、不在のファイル・データブロックをそのデータブロックを有するサーバ110に要求する。図1Cの例では、不在のファイル・データブロックが破線ボックス118cとして図示されている。上記の例では、サーバ・ファイルハンドル「FILE1 ON SERVER1」を用いてその要求が送信される。

【0040】

NFSフィルエンジン166がサーバ110からファイル・データブロックをフェッチ（197）すると、そのデータブロックは、リモート・フィル層168、ファイルシステム169およびストレージデバイス・マネージャ170によってストレージ・ディスク171に入れられる。そしてバッファ・ツリー191のブロック番号が更新され、ファイル・データブロックが割り当てられる。次に、それらのファイル・データブロックが要求元クライアント105に送信され、要求されたファイルがそのクライアントに与えられる。上記の方法は、ディレクトリ要求の処理にも使用することができる。

10

【0041】

また、図1Cに示す実施形態によれば、バッファ・ツリー内のファイル・データブロックを1つのファイル内にまばらにキャッシュすることができる。この方法は、ファイルの途中における丸め処理や書込み処理のような、部分的なファイル処理を実施することが可能になるという有利がある。

20

【0042】

通常、クライアント105がNFSプロトコルで要求を送信する場合、各データブロックについて個別の要求が作成される。また、CIF Sプロトコルの場合にも、各データブロックについて個別の要求が作成される。例えば、第1のクライアント105（例えば、図1のクライアント105（1））がデータブロック188aおよび188bを求めるデータ要求を送信するものと仮定する場合、それらのデータブロックは両方とも図1Cの例におけるプロキシ・キャッシュ115にキャッシュされている。クライアントからのデータ要求の数は変更することもできる点に注意して欲しい。データブロック188aおよび188bがプロキシ・キャッシュ115にキャッシュされているので、プロキシ・キャッシュ115がクライアント要求を受信するとキャッシュ・ヒット状態が発生し、要求されたデータブロック188aおよび188bが、プロキシ・キャッシュ115から要求元クライアント105（1）に送信される。

30

【0043】

キャッシュ・ミス状態の例として、第2のクライアント105（例えば、図1のクライアント105（2）など）が、データブロック188a、188bおよび188cのデータを求める要求を送信する場合を仮定する。図1Cの例では、プロキシ・キャッシュ115がクライアント要求を受信したとき、データブロック188aおよび188bはプロキシ・キャッシュ115にキャッシュされているが、データブロック188cはプロキシ・キャッシュ115内に存在しない。プロキシ・キャッシュ115はこのキャッシュ・ミス状態に回答して、上記と同様にサーバ110からデータブロック118cをフェッチする。フェッチされたデータブロック188cは、プロキシ・キャッシュ115に割り当てられた後、プロキシ・キャッシュ115から要求元クライアント105（2）に送信される。

40

【0044】

さらに他の例として、他のクライアント105（例えば、図1のクライアント105（n））がデータブロック188bおよび188cを求めるデータ要求を送信する場合を仮定する。データブロック188cはサーバ110からプロキシ・キャッシュ115に前もってフェッチされ、プロキシ・キャッシュ115に割り当てられているので、プロキシ・キャッシュ115がクライアント要求を受信したときに、プロキシ・キャッシュ115はデータブロック188bおよび188cを有している。従って、キャッシュ・ヒット状態

50

が発生し、要求されたデータブロック 188b および 188c が、プロキシ・キャッシュ 115 から要求元クライアント 105 (1) に送信される。クライアント 105 がプロキシ・キャッシュ 115 内に存在しないデータブロック 188 を要求した場合、プロキシ・キャッシュ 115 は、その存在しないデータブロック 188 をサーバ 110 からフェッチし、キャッシュしてから、そのフェッチしたデータ 188 を要求元クライアント 105 に送信する。

【0045】

このように、プロキシ・キャッシュ 115 は、クライアント要求に応答して部分オブジェクトの複製を行なう。部分オブジェクトは、ファイルの 1 以上のデータブロックとして定義される。上記の例では、クライアント要求に応答してデータブロック 188c をプロキシ・キャッシュ 115 に複製した後、それを他の要求元クライアント 150 から使用できるようにしていた。以下で説明するように、適当なデータ置換ポリシーを用いて、(プロキシ・キャッシュ 115 から) データブロック 188、または複数のデータブロック 188 によって定義されるファイルを消去することができる。これに対し、従来の複製技術は、全てのファイルを複製している(例えば、種々のミラーリング技術は、全てのボリュームまたは全てのファイルセットを複製している)。

10

【0046】

置換ポリシー

プロキシ・キャッシュ 114 においてアクティブデータセット 140 内のデータファイルを更新、維持、または消去するための置換ポリシーには、様々なものがある。図 2 は、アクティブデータセット 140 内にキャッシュされたデータの重要度を、そのキャッシュされたデータに対するクライアント(複数の場合もあり)からのアクセスすなわち要求の回数に基づいて増加させる方法を示すグラフ 200 である。従ってプロキシ・キャッシュ 115 は、アクティブデータセット 140 内にキャッシュされたデータに対し、そのキャッシュされたデータに対するクライアント 105 からのアクセスまたは要求の回数に基づいて「重要度」値を動的に割り当て、調節するように構成される。プロキシ・キャッシュ 115 は、アクセス数すなわち要求回数が増えるのに従って、キャッシュされたデータの重要度を増加させる。例えば、プロキシ・キャッシュ 115 は、アクティブデータセット 140 内にキャッシュされたデータが重要度値 Y (図 2 のグラフに点 205 で示す) に達した場合、プロキシ・キャッシュ 115 のストレージユニットから、そのキャッシュ・データが消去されないようにする。

20

30

【0047】

図 3 は、本発明の一実施形態による、アクティブデータセット 140 を管理する他の方法 300 を示すブロック図である。上記のように、アクティブデータセット 140 は、先入れ先出し(FIFO)法 305 などの適当な方法により、置換、またはプロキシ・キャッシュ 115 から消去することができる。FIFO 法 305 では、プロキシ・キャッシュ 115 のストレージユニットから消去する際に、アクティブデータセット 140 を循環させる。最長時間未使用(LRU)法などの他の適当な置換ポリシーを使用することもできる。LRU による置換方法は、一般にデータベース管理システムで使用される方法で、最も長時間使用されなかったブロックを最初に置換すべきブロックとする方法である。

40

【0048】

代替または追加として、アクティブデータセット 150 内のファイルは、ユーザがクライアント 105 からプロキシ・キャッシュ 115 へロックコマンドを送信することにより、ロック(310)することもできる。ロックコマンドは、プロキシ・キャッシュからそのファイルが消去されないようにするためのものである。ユーザは、例えば、サーバ 110 までのリンク接続が壊れているかもしれないとか、サーバ 110 が故障しているかも知れないといったことが心配である場合に、ロックコマンドを送信する場合がある。代替として、オブジェクト(例えば、ファイルや、ファイルのデータブロックなど)の属性(メタデータ)は、そのオブジェクトがプロキシ・キャッシュ内にどのくらいの期間キャッシュされているのかを示すことができる。

50

【0049】

代替または追加として、ユーザは、ヒントまたはインジケータをプロキシ・キャッシュ115に与え(315)、特定のキャッシュ・ファイルを重要なものとして指定することもできる。その結果、プロキシ・キャッシュ115のストレージユニット(複数の場合もあり)から指定された重要ファイルが消去されることが、インジケータによって防止される。

【0050】

プロキシ・キャッシュ115を使用すると、データを様々な遠くの場所に分散および複製することができ、データ複製のための従来のミラーリング技術の使用を回避できるという利点がある。従来のミラーリング技術を使用するためには、所定の時刻または所定の時間間隔(例えば毎晩)で全てのデータセットをミラーリングしなければならない、全てのデータセットを入れるディスク空間が必要となる。これに対し、プロキシ・キャッシュ115は、アクティブデータセット140内のデータをオン・デマンドで複製する。アクティブデータセット140は、従来のデータ複製方法に関する大きなディスクの必要性を無くすことができるという利点を有する。さらに、アクティブデータセット140はキャッシュ・ミスに回答して自動的に更新または消去されるので、従来のデータ複製方法に関する特殊なソフトウェアの必要性も無くすことができる。

【0051】

マウントポイントの統合とファイルハンドルの書換え/変形の方法

NFSネットワーク・ファイルシステム・プロトコルによって処理を実行するために、クライアントは、NFS要求を下記の(1)~(3)とともにNFSサーバに送信する。

- (1) 処理対象を指定するためのNFSファイルハンドル
- (2) 処理(検索、読出し、書込み、パーミッション変更など)
- (3) 要求を送信する代表となるユーザ

NFSクライアントがリモートのファイルシステムに最初にアクセスしようとする場合、クライアントはまず、ボリューム(上で述べたように、「ボリューム」という用語と「ファイルシステム」という用語は、同じ意味で使用される)のエントリ・ポイントのファイルハンドルであるルート・ファイルハンドルを入手しなければならない。その目的のために、クライアントホストは、マウント要求をサーバのマウントデーモンに送信する。ここで、マウント要求とは、ルート・ファイルハンドルにアクセスするためのプロトコルの一部であり、デーモンとは、コンピュータシステムが受信する可能性のある定期的なサービス要求を処理するために、絶え間なく動作し続け、存在するプログラムである。デーモンプログラムは、必要に応じてその要求を他のプログラム(またはプロセス)に転送する。サーバのマウントデーモンは、要求されたファイルシステムにアクセスするためのパーミッションをクライアントが有していることを確認する。マウントデーモンは、アクセスを許可する場合、ファイル(ディレクトリ)ハンドルをそのNFSクライアントに返送する。通常、各ファイルハンドルは、32バイトの不可視の識別データである。ファイル名が変更されても、そのファイルハンドルは、リネーム後のファイルのファイルハンドルとしてそのまま維持される。従って、上記のマウント処理は、ルート・ファイルハンドルにアクセスするためのプロトコルである。

【0052】

2つの異なるサーバが同じファイルハンドルを使用することもありうるので、NFSのファイルハンドルは、世界的に/いかなる場合にも、重複しないものではない。従来は、この条件が問題となることはなかった。その理由は、各サーバにどのファイルハンドルが関連しているのかを、クライアントが常に管理しているからである。プロキシ・キャッシュ115を使用して複数のマウントポイントを統合する場合、プロキシ・キャッシュ115は、多数の異なるサーバから得たボリュームをエクスポートする場合がある。それらのファイルハンドルは必ずしも重複していないとは限らないので、プロキシ・キャッシュ115は、どのファイルハンドルがどのサーバから得られたものであるのかを判断できない場合があり、第1のサーバにおける第1のファイルと第2のサーバにおける第2のファイ

10

20

30

40

50

ルが同じファイルハンドルを有することもあり、その状態は衝突問題を引き起こす可能性がある。この問題を悪化させる原因は、ファイルハンドルが不可視なものとして定義されていることにある。ここで、不可視とは、ファイルハンドルの内容を知ることができず、ファイルハンドルの名前だけしか得られないことを意味している。

【 0 0 5 3 】

本発明の一実施形態では、ファイルハンドルをクライアント 1 0 5 に送信する前に、ファイルハンドルを変形することにより、プロキシ・キャッシュ 1 1 5 がファイルハンドルを解釈してファイルハンドルの宛先サーバ 1 0 5 を判断できるようにする。従ってクライアント 1 0 5 は、データの入手元であるサーバを明確に知らなくても、プロキシ・キャッシュ 1 1 5 を介してデータを得ることができる。

10

【 0 0 5 4 】

図 4 に示すように、ファイルハンドル 4 0 0 は、F S _{i d} (ファイルシステム識別子) フィールド 4 0 5、ファイル識別子 (I D) フィールド 4 1 0、生成番号フィールド 4 1 5、および、その他データ用フィールド 4 2 0 などの種々のフィールドを有する。通常、F S _{i d} フィールド 4 0 5 は約 4 ~ 8 バイトである。

【 0 0 5 5 】

本発明の一実施形態によれば、サーバ側ファイルハンドルの仮想化が可能になる。クライアント 1 0 5 とサーバ 1 1 0 との間にプロキシ・キャッシュ 1 1 5 を導入すると、ファイルハンドルを変形 (変換) し、世界的に重複しない名前空間を構成することが可能になる。クライアントに返送されるファイルハンドルは、それらが種々の入手元サーバまたは種々のエクスポートオプションセットにマッピングされること、すなわち、複数のマウントポイントを交差させるものであることを示すものに変更される。クライアントに送信されるファイルハンドルと、入手元サーバ上のオブジェクトのアクセスに使用されるファイルハンドルとの間に間接層を追加することにより、クライアント 1 0 5 に何も影響を与えることなく、入手元サーバ 1 1 0 に対する変更を行なうことができる。

20

【 0 0 5 6 】

図 5 および図 6 は、本発明の一実施形態による、マウントポイントを統合する方法およびファイルハンドルを書換え、すなわち変換する方法を示す図である。図 5 に示すように、各サーバ 1 1 0 は通常、ファイルシステムに相当するボリューム (例えば、ボリューム / V o l 1、ボリューム / V o l 2、ボリューム / V o l 3 など) を含む。つまり、サーバ 1 1 0 (1) はボリューム 5 0 5 ~ 5 1 5 を含み、サーバ 1 1 0 (2) はボリューム 5 2 0 ~ 5 3 0 を含む場合がある。

30

【 0 0 5 7 】

本発明の一実施形態では、図 6 に示すように、特定のサーバ 1 1 0 の各ボリュームについて F S _{i d} 値のマッピング 6 0 0 を作成し、そのマッピング 6 0 0 をプロキシ・キャッシュ 1 1 5 に格納する。あらゆる F S _{i d} 値を変換することで、キャッシュ・ミスが発生したときに、あるファイルがどのサーバに格納されているかを判定することが可能になるという利点がある。

【 0 0 5 8 】

変換テーブル 6 0 5 は、ハッシュ値 N (例えば、N 1、N 2、N 3、・・・など) を有し、F S _{i d} テーブル 6 1 0 は、サーバ 1 1 0 の各ボリュームについて F S _{i d} 値を有する。ファイルハンドル中に特定の F S _{i d} 値を有するサーバ 1 1 0 (例えば、図 6 の例の場合、F S _{i d} = 値 1 を有するサーバ 1 1 0 (1) のボリューム 5 0 5 など) からトラフィック 6 1 5 が受信されると、その F S _{i d} = 値 1 は、プロキシ・キャッシュ 1 1 5 のテーブル 6 0 5 の F S _{i d} 値 N 1 に変換される。サーバ 1 1 0 から受信された他のファイルハンドル中の他の F S _{i d} 値は、テーブル 6 1 0 に格納され、テーブル 6 0 5 内の様々な値に変換、すなわち変形される。従って、プロキシ・キャッシュ 1 1 5 は、クライアント 1 0 5 からファイルハンドル 6 2 0 を受信したとき (クライアントのデータ要求の際) にキャッシュ・ミスが発生すると、テーブル 6 1 0 に N の値および F S _{i d} の値を有するマッピング 6 0 0 を用いて、そのファイルハンドル 6 2 0 を送信すべき特定のサーバ 1 1 0

40

50

およびファイル/ボリュームを判定する。プロキシ・キャッシュ115は、ファイルハンドル620を宛先サーバ110に送信する前に、ファイルハンドル620の変換後のFS_{id}値Nをテーブル610中の適当なFS_{id}値に変換する。

【0059】

テーブル600における値Nは、ボリュームおよび/または他の識別子を有するサーバ110のアドレスから何らかのハッシュを利用して計算される。各値Nに関するハッシュは重複しないので、上記のような名前衝突問題は起こらない。値Nは通常、8バイトの数値であり、サーバ名、ボリューム名およびファイル名(例えば、ファイラ1/ボリューム0/ファイル0など)によって決まる設定値を有する。

【0060】

あるいは、値Nは、上で述べたように、サーバのインターネット・プロトコル(IP)アドレスおよびサーバのファイルハンドル値からMD5アルゴリズムを用いて計算することもできる。ファイルハンドル値は種々のサーバ間で重複しない値であるとは限らないので、サーバIPアドレスを用いて作成すれば、ファイルに関する各Nの値をサーバ間で異なる値にすることができる。

【0061】

図7に示すように、マッピング600は様々なプロキシ・キャッシュ(例えば、プロキシ・キャッシュ115と705)に複製することができる。そのため、あるプロキシ・キャッシュが故障した場合でも、他のプロキシ・キャッシュを利用することができる。例えば、装置故障によりクライアント105がプロキシ・キャッシュ115にアクセスできなくなった場合、プロキシ・キャッシュ705へのフェイルオーバー710を実施することができ、その結果、クライアント105はファイルハンドル400の書換え用のマッピング600を用いてプロキシ・キャッシュ705にアクセスできるようになる。

【0062】

マッピング600は複数のプロキシ・キャッシュに複製することができるので、代替または追加として、クライアント105がボリュームをアンマウント・リマウントしたり、他の変更をクライアント105に施したりしなくても、プロキシ・キャッシュ(例えば、プロキシ・キャッシュ115など)は、新たなプロキシ・キャッシュまたは他のプロキシ・キャッシュ(例えば、プロキシ・キャッシュ705など)に交換することができる。

【0063】

世界的に重複しない仮想名前空間を作成して使用方法

本発明の一実施形態によれば、上記のようなマウントポイント間の統合が可能になるだけでなく、世界的に重複しない名前空間の作成も可能になる。従来は、NFSのマウントポイントをネストさせること(すなわち、クライアントがNFSボリュームをマウントするときに、そのボリュームに他のマウントポイントを含めること)が出来なかった。この制約は、ネットワークのユニフォーム・ビューの作成を非常に難しいものにしていた。

【0064】

本発明の一実施形態によるファイルハンドルの仮想化を使用すると、管理者は、NFSによってエクスポートされたボリューム内の任意のディレクトリをマウントポイント(ファイルシステムに対するアクセスポイント)として設定することが可能になる。これは、クライアント105が1つのボリュームをマウントするだけで、誰かがそのボリュームにアクセスしたときに、そのクライアント要求がプロキシ・キャッシュによって特定の適当なボリュームに適当に宛先変更されることを意味する。従って、プロキシ・キャッシュ115は、単一のマウントポイントを用いて、各サーバ110(または新たに追加されたサーバ)についてあらゆるクライアント105に通知することができる。オブジェクトをリネームすることにより、ファイルの仮想グローバル・ビュー(すなわち、仮想グローバル名前空間)が作成される。その結果、各クライアント105は、ファイルの仮想グローバル・ビューを用いて、各サーバ110に関する情報および各サーバ110の各ファイルに関する情報を得ることができる。仮想グローバル名前空間を使用すると、クライアント105の管理が簡単になる。

10

20

30

40

50

【0065】

図8Aに示すように、(プロキシ・キャッシュ内の)ボリューム名805は、複数のクライアント105からアクセスすることができる。ボリューム名の例は、例えば「/global」である。ボリューム名805は、本発明の一実施形態による仮想名前空間を可能にする。クライアント105は、ボリューム名805、並びにサーバ110(1)、110(2)、110(m)にそれぞれ関連するビュー・フォルダ810、815、820mをマウントすることができる。各フォルダ810、815、820mには、特定サーバに割り当てられたボリューム名805が格納される。従ってボリューム名805は、サーバ110(1)の/vol505、サーバ110(1)の/vol510、サーバ110(1)の/vol515、サーバ110(2)の/vol520、サーバ110(2)の/vol525、サーバ110(2)の/vol530、およびサーバ110の他の既知のボリュームを指し示す。従って、フォルダ805~820mは、ボリューム名805によってマウントポイントに仮想的にマッピングされる。クライアントは、1つのマウントポイント(例えば、/globalなど)を知っているだけで、適当なサーバ110にアクセスすることができる。マッピング600(図6を参照)を用いると、NFSプロトコルにおけるファイルハンドルの書換えにより、ファイルハンドル400をサーバ100のボリューム内の仮想アクセスポイントとして使用することが可能になる。上記のように、マウント処理は、サーバと通信し、ファイルシステムのルート・ファイルハンドルを得ることを含む。このルート・ファイルハンドルは、後に検索リモート・プロシージャ・コール(RPC)に渡され、リモートファイルシステムのディレクトリ階層中の他のファイルハンドルを見付けるのに使用される。

【0066】

単一マウントポイントによって得られる利点は、クライアント105があらゆる単一マウントポイントを知っている必要はない点にある。従来の方法では、各クライアントが通常、フォルダ内の既知のリモートファイルシステムが全てリストされたfstabファイル(ファイルシステム・テーブルファイル)を有している(すなわち、fstabは、クライアントがドライバ、サーバ名およびボリューム名にアクセスできる場合の、クライアントにおけるローカル・ディレクトリまたはローカル・フォルダのリストである)。サーバ110に変化が生じ、その変化が名前空間に影響を与えた場合、各クライアント105を再構成し、クライアントがサーバにアクセスできるようにしなければならない。そのため、従来の方法におけるマウントポイント管理作業は、複雑で時間を要するものとなっている。

【0067】

本発明の一実施形態の場合、各クライアントはディレクトリ805をマウントするだけでよく、この1つに統合されたマウントポイントによって、マウントポイントの管理を単純化している。特に、このマウントポイントは、プロキシ・キャッシュ115に対してアンマウントおよびマウントされる。そのため、ボリューム805をクライアント105に対してアンマウントおよびリマウントする必要はない。従って、ディレクトリ805はファイルシステムテーブルファイルの機能を提供し、クライアントがドライバおよびサーバにアクセスできるようにする。ディレクトリ805は、リモート・ファイルシステムへのアクセスポイント、および、リモートサーバ(複数の場合もあり)110に格納されているような属性のリストを提供する。

【0068】

複数のマウントポイントを作成し、属性集合をマウントポイントに関連付ける方法

同一オブジェクトについて、ファイルハンドルは様々な方法で書き換えることができる。一実施形態において、ファイルハンドルは、プロキシ・キャッシュ上の関連ローカルマウントポイントの属性を利用して様々な方法で書き換えることができる。

【0069】

次に、図8Bを参照する。マウントポイントは、ローカルボリューム名として定義される。例えば、ローカルボリューム名「user/read-only/employee

1」(ブロック850を参照)とローカルボリューム名「user/employee1」(ブロック855を参照)はそれぞれ、2つの異なるファイルハンドルを生成する。その理由は、各ローカルボリューム名についてマウントポイントが異なるからである。具体的には、「user/read-only/employee1」850のマウントポイントは「user/read-only」である一方、「user/employee1」855のマウントポイントは「user」である。

【0070】

ローカルボリューム名は、関連する属性またはメタデータを有することができる。例えば、アクセス制限、格納するための資格、またはキャッシュ機能、セキュリティ属性、認証、プロキシ・キャッシュ装置の認証レベル、サーバの読出し専用属性または書込み専用属性などである。データを共有する他のボリュームに関連するマウントポイントにアクセスすることにより、そのデータに関連する属性を変更することができる。つまり、個々の属性は特定のマウントポイントまたはアクセスポイントに関連付けることができ、それをハッシュすると、重複しない識別子が形成される。一実施形態において、この重複しない識別子は、ローカルボリューム名、ファイルシステムIDおよびサーバIPアドレスをMD5のような適当なハッシュアルゴリズムを用いてハッシュしたものである。図8Bの例では、属性(複数の場合もあり)856がマウントポイント850に関連付けられ、属性(複数の場合もあり)857がマウントポイント855に関連付けられている。

10

【0071】

なお、プロキシ・キャッシュ115のディスク上の同じ物理オブジェクトを異なるマウントポイントからアクセスすることも可能である。

20

【0072】

キャッシュ・データの整合性をとる方法

最初に、プロキシ・キャッシュがNFSクライアントに対して保証するデータ整合性モデルを定義する。このモデルは、NFSサーバによって保証されるものと同じである。次に、サーバおよびプロキシ・キャッシュがどのようにしてこのNFSデータ整合性モデルを保証するのかについて説明する。次に、入手元サーバにアクセスし、キャッシュとディレクトリの両方を通してデータを変更しているクライアント群に対し、このデータ整合性モデルがどのように作用するかについて説明する。最後に、キャッシュ最大寿命(すなわち、入手元サーバを用いて新しさが確認される前にヒットが満たされる際の最大量の時間)を使用する方法について説明し、キャッシュ最大寿命がデータの新鮮さおよび整合性にどのように影響を与えるかについて説明する。以下に記載する例はNFSプロトコルに関するものになっているが、キャッシュ・データの整合性をとるためのこの方法は、他のプロトコルでも使用することができる。

30

【0073】

NFSクライアントデータ整合性

クライアントの視点から見ると、NFS整合性モデルは、番号順のファイル・トランザクションを保証するものである。あるトランザクションが完了し、特定バージョンのファイルが返された後、そのファイルに対して開始されるトランザクションはいずれも、少なくとも新しいバージョン、すなわちもっと新しいバージョンを返すことが保証される。我々はこれを連続整合性と呼んでいる。NFSにおける整合性の理解が難しい理由は、例えば、種々の時間にわたるトランザクションがパイプライン化され、信頼性のないトランスポートに依存することがあるからである。従って、下記のように少数のクライアントの定義から始めることが有用である。

40

【0074】

トランザクションの開始

NFSトランザクションの開始は、クライアントが最初の要求RPC(リモート・プロシージャ・コール)を送信した時刻として定義される。以後のDUP要求(複製要求)がトランザクションの開始をリセットすることはない。DUP要求は、NFSプロトコルによって以前に行なわれた要求の複製である。ファイルサーバ上のDUPキャッシュは、ク

50

クライアントからの要求に対する応答だけでなく、トランザクション識別子 (X I D) も格納している。各 N F S 要求は、重複しない X I D を有する (要求が重複しない X I D を有するのに対し、応答は要求に対応する X I D を有する)。複製要求は、もとの要求と同じ X I D を有する。

【 0 0 7 5 】

トランザクションの終了

N F S トランザクションの終了は、クライアントが、発生しうる複数の D U P 応答のうち最初のものを受信したときに発生する。

【 0 0 7 6 】

連続トランザクション

一方のトランザクションが完了する前に他方のトランザクションが開始された場合、それら 2 つのトランザクションは連続であるものと定義される。

【 0 0 7 7 】

ファイル

N F S サーバ上では、ファイルは、そのファイルのアクセスに使用されるプロトコル・バージョンやエクスポートパスに応じて、種々のファイルハンドルによって参照される。プロキシ・キャッシュの一実施形態において、ファイルはファイルハンドルによって一意に識別される。これは、異なるエクスポートパスまたは異なるプロトコル・バージョンを用いて複数のマウントポイントにインポートされた実際には同一のファイルが、複数の別個のファイルとみなされることを意味している。

【 0 0 7 8 】

整合性

サーバの視点から見ると、サーバは通常、クライアント側でトランザクションが何時開始され、何時終了するのかを正確には知らない。サーバは応答を送信した後、それがクライアントに直ぐに受信されることを想定している。サーバが要求を受信した場合、サーバは、以前に送信した応答がすべてクライアントに受信された後で、その要求をクライアントが開始したものと想定する。サーバの視点から見れば、トランザクションは最初の要求を受信した時に開始され、最初の応答を返送した時に終了する。

【 0 0 7 9 】

サーバ上で実行しなければならない変更処理は、一度だけであることが望ましい。例えば、 R E M O V E 処理について考えてみる。 R E M O V E 処理が最初にサーバに受信されたときは、ファイルシステムから適当なファイルが消去され、その記録が D U P キャッシュに作成される。以後の複製要求に対する応答はすべて、この D U P キャッシュを用いて行なわれる。これによって、以前の応答の中にクライアントに届くまでの間に失われたものがあっても、複製要求の失敗を防止することができる。

【 0 0 8 0 】

N F S プロキシ・キャッシュにおける整合性

プロキシ・キャッシュがクライアントに応答を送信した場合、プロキシ・キャッシュは、それ以降に受信した全ての要求に対する応答が、整合性のとれた返答になることを保証しなければならない。これは、特定バージョンのファイルをクライアントに教えた後、そのファイルに関する以後のトランザクションは全て、少なくともその特定バージョンと同程度に新しいバージョンを返さなければならないことを意味している。

【 0 0 8 1 】

プロキシ・キャッシュの視点からすれば、これは次の 2 つを保証することを意味する。

(1) プロキシ・キャッシュは、応答をクライアントに送信するときに、以後のキャッシュ・ヒットに対して整合性のとれた応答を返すように準備することが望ましい。これは、従来は、キャッシュ・ミス時に応答をクライアントに送信する前に、データが満杯であることを確認することによって行なわれている。

(2) プロキシ・キャッシュにデータを入れる場合、プロキシ・キャッシュは、クライアントに既に返送されたそれよりも新しいバージョンのファイルが無いことが確認できた

10

20

30

40

50

場合にのみ、そのデータを以後のキャッシュ・ヒット時に使用することができる。我々は、これを整合性充填と呼んでいる。整合性充填については以下の実施形態で説明する。

【0082】

キャッシュ・リブートの扱い

プロキシ・キャッシュが起動すると、プロキシ・キャッシュはメモリ内オブジェクトキャッシュを再構築しながらプロキシ専用モードで動作する。一実施形態において、充填はロギングされず、不適切なシャットダウンは、キャッシュを不整合状態にすることがある。そのため、プロキシ・キャッシュは起動時に、全てのキャッシュ・オブジェクトを、キャッシュ・ヒット時に返送する前に確認が必要なものとしてマークする。プロキシ・キャッシュは、クライアント要求に回答してオブジェクトをクライアントに渡す前に、そのオブジェクトの属性をチェックすることにより、そのオブジェクトが最新のバージョンであることをサーバを用いて確認する。

10

【0083】

リポート時には、プロキシ・キャッシュは全てのトランザクションを中止し、トランザクションがクライアントによって再起動される。サーバから配送中のキャッシュ充填要求はプロキシ・キャッシュによって捨てられる（破棄される）。なぜなら、それらのキャッシュ充填要求のXID（トランザクション識別子）は通常、解釈できないものになるからである（リポートの際には通常、後続の要求のXIDについて新たなXIDマッピングがプロキシ・キャッシュ内に作成されるので、新たなXIDマッピングはリポート発生前のXIDを解釈することができない）。リポートの後、到来するクライアント要求はすべて、新たなトランザクションの最初の要求として扱われる。起動時に全てのトランザクションを再起動させれば、プロキシ・キャッシュは、すべてのアクティブ・トランザクションについて開始時刻を記録することができる。これによって、整合性充填に関する本発明の一実施形態が可能になる。

20

【0084】

整合性キャッシュ充填

キャッシュ充填には次の2つがある。

(1) 読み出し専用処理によって返送されたデータの格納

(2) キャッシュ・データに対する変更処理の再現

変更処理は通常、プロキシ・キャッシュで常に見られるものではないので、あるいは、それらを再現することは非常に難しいので、プロキシ・キャッシュからファイルを投げなければならぬ場合がしばしばある。「投げる」とは、特定タイプの充填を簡略化したものである。また、性能上の理由から、データの充填が有用な場合もあり、そのデータは新しい場合もあれば、キャッシュファイルを不整合状態にする場合もある（NFS応答によっては、参照される全てのファイルについてファイル属性を持たないものもある）。本発明の一実施形態では、データを充填することによってこの状態を処理する。ただし、ファイルは、後続のヒット時にキャッシュから返送される前に確認が必要なものとしてマークされる。

30

【0085】

整合性のとれたデータがクライアントに常に送信されることを保証するためには、全てのファイルを次の2つのうちの一方に分類しなければならない。

40

(1) 整合性があるもの

(2) 不整合の可能性のあるもの

不整合の可能性のある充填は、以前にクライアントに送信されたデータに比べて、新しいデータを有する場合がある。NFSバージョン2では、応答が新しいデータを有しているか否かを常に知ることができるとは限らない。この種の応答は、不整合な充填の候補となる。応答に関連するファイル（複数の場合もあり）に変更があった時にサーバから再要求された応答も、候補となる。プロキシ・キャッシュのアプローチは通常、堅実なものではないので、確実なことが言えない場合は「新しい」と推測する。不整合な充填があると、プロキシ・キャッシュは、後続のヒット時に問題のファイルを提供する前に

50

、その新しさを入手元サーバを用いて確認しなければならない。これは、直後に簡単なキャッシュ・ヒットが続く可能性がある整合性のある充填とは対照的である。

【0086】

一実施形態では、全てのトランザクションが起動時に再起動され、キャッシュされた内容はすべて、確認が必要なものとしてマークされる。これによって、全てのキャッシュデータが、整合性のあるものとして開始されることが保証される。ファイルFに対する不整合な充填は、Fに対する最後に記録された変更よりも前にキャッシュによって開始されたトランザクションの結果得られるものとして定義される。最後に記録された変更は、変更の可能性を示すサーバ応答がプロキシ・キャッシュによって受信された時刻として定義される。これは、その変更がプロキシ・キャッシュのクライアントから見れるようになった時である。

10

【0087】

ファイルに対する最後に記録された変更を追跡する場合、プロキシ・キャッシュは非常に堅実なものであることが望ましい。プロキシ・キャッシュは、サーバからの応答を処理するとき、例えば下記の事項のうちのいずれかが真であれば、トランザクションに関連するファイルに対する最後に記録された変更を更新する。

- ・ 応答が、ファイルに対する変更があったことを知っている処理に関するものである。
- ・ 応答が、ファイルの処理後属性を含まない（処理後属性とは、ファイルに対して処理を実施した後のファイルの属性である）。
- ・ 応答が、ファイルの処理後属性を含み、それらの属性が、現在キャッシュされているものとは異なる。
- ・ ファイルの属性が何もキャッシュされていない。
- ・ キャッシュが再構築段階にある。

20

【0088】

各ファイルについて最後に記録された変更を追跡する本方法は、ハッシュテーブル（下記の説明を参照）を用いて重複しないファイル名（例えば、MD5アルゴリズムによって得られるものなど）を定めることによる。ファイルは通常、LRUに返還要求される。トランザクションがこのLRUにおける最も古い要素よりも古い場合、そのトランザクションに関する最後に記録された変更は返還要求された可能性があるものと仮定され、最後に記録された変更は不整合な充填として扱われる。

30

【0089】

分散NFS整合性

ファイルシステムを求める全てのクライアントが特定のプロキシ・キャッシュを通過するとすれば、それらのクライアントは、データの最新のコピーを得ることが保証される。その理由は、変更処理はすべて、プロキシ・キャッシュで実施されるか、または、新たなキャッシュデータが投げられるかのいずれかとなるからである。具体的には、全ての変更処理がプロキシ・キャッシュを通過するとすれば、たとえ読み出し専用処理をルーティングしたとしても、このことが言える。

【0090】

クライアントが様々な経路を通じてファイルシステムにアクセスする場合（一部がプロキシ・キャッシュを通り、一部が入手元サーバに直接アクセスする場合など）、複雑さが生じる。プロキシ・キャッシュは、変更処理を入手元サーバに直接渡した後、新しいデータをクライアントに提供することができる。プロキシ・キャッシュのクライアントが新しいデータを得ることができ、且つ、整合性のとれたデータを常に受け取れることは重要である。クライアントによって様々な新しさの許容範囲が認められるためには、プロキシ・キャッシュは、様々なデータセットについて最大寿命を指定することが出来なければならない。この最大寿命は、入手元サーバを用いて新しさが確認される前にヒットが満たされる際の最大量の時間である。

40

【0091】

図8Cは、プロキシ・キャッシュがサーバに対して前回、要求R0を行い、その要求R

50

0 が時刻 A 0 でサーバに受信された場合の図である。プロキシ・キャッシュが次の要求 R 1 をサーバに対して行くと、その要求は時刻 A 1 でサーバに受信される。ただし、要求 R 0 に応答するデータ D 0 は、要求 R 1 に応答するデータ D 1 をプロキシ・キャッシュが受信した後で、プロキシ・キャッシュに受信される。従って、後続の受信データ D 0 は、ファイル D の最新バージョンではないかもしれない。

【0092】

図 8 D は、プロキシ・キャッシュがサーバに対して前回、要求 R 0 を行い、その要求 R 0 が時刻 A 0 でサーバに受信された場合の図である。プロキシ・キャッシュが次の要求 R 1 をサーバに対して行なうと、その要求は時刻 A 1 でサーバに受信される。要求 R 0 に応答するデータ D 0 は、要求 R 1 に応答するデータ D 1 をプロキシ・キャッシュが受信する前に、プロキシ・キャッシュに受信される。

10

【0093】

図 8 E に示すように、本発明の一実施形態では、ハッシュテーブル 8 7 0 をプロキシ・キャッシュに格納する。ハッシュテーブル 8 7 0 はエントリ 8 7 2 を有する。各エントリは、オブジェクト 8 7 3 (例えば、ファイルやデータブロックなど) および最後に記録された変更値 8 7 4 を含む。例えば、返却データ D 0 は、その最後に記録された変更を示す属性 L R C 0 を有する場合がある(そのオブジェクトが最後に変更されたものである場合)。同様に、返却データ D 1 は、その最後に記録された変更を示す属性 L R C 1 を有する場合がある。本発明の一実施形態において、データ D 0 の L R C 0 値が、所定の時間(例えば、30 秒など)内に要求 R 1 が成された時間よりも大きい場合、プロキシ・キャッシュは D 1 の属性を確かめるための要求を作成し、D 1 がデータ D の最新バージョンであるか否かを確認する。一実施形態において、プロキシ・キャッシュ 1 1 5 のフィルエンジン 1 6 6 は、例えば「GETATTR」コマンドをサーバに送信し、データ D の属性を得る。この方法によれば、プロキシ・キャッシュに格納されたデータのデータ整合性を達成することが可能になる。

20

【0094】

プリフェッチによる属性の事前確認

HTTP プロトコルの場合、サーバは通常、クライアントがサーバにおけるファイル更新をチェックしなければならない時を示すインジケータをクライアントに与える。これに対し、NFS プロトコルでは、サーバは、サーバファイルの更新をクライアントに通知することができない。以下で説明する種々の事前確認方法は、プロキシ・キャッシュがサーバに格納されたデータの属性を確認できるようにするためのものである。

30

【0095】

プロキシ・キャッシュ 1 1 5 を用いると、サーバ 1 1 0 とクライアント 1 0 5 が離れている場合に(例えば、サーバ 1 1 0 とクライアント 1 0 5 が異なる都市にある場合に)、WAN 1 2 5 を介してデータを統合することができる。プロキシ・キャッシュ 1 1 5 は、クライアント 1 0 5 がサーバ 1 1 0 に格納されているデータを要求する際の遅延を減らすことができるので、サーバ 1 1 0 はクライアント 1 0 5 にとってローカルであるかのようにになる。遅延を減らすために、本発明の実施形態は、NFS のようなオープン・プロトコルを WAN 1 2 5 を介して確実に伝送できるようにする解決策を提供する。NFS を使用するクライアントは通常、フォルダ内の各ファイルに関する情報を順番に要求する。以下で説明するような本発明の一実施形態によるプリフェッチ方法は、WAN を介した NFS の遅延を減らすことができるという利点がある。図 9 に示すように、クライアント 1 0 5 がデータの要求 9 0 5 を送信し、キャッシュ・ミスが発生した場合、プロキシ・キャッシュ 1 1 5 は、サーバ 9 1 0 内のデータの属性をプリフェッチ 9 1 0 し、サーバはその属性を返す 9 1 5。そのため、プロキシ・キャッシュ 1 1 5 は、データのいかなる変化も判定することができる。その後、プロキシ・キャッシュ 1 1 5 は、要求されたデータをクライアントに送信することができる。

40

【0096】

本発明の一実施形態では、サーバ(複数の場合もあり) 1 1 0 内のフォルダ、サブフォ

50

ルダ、および/または、ファイルの属性をチェックすることにより、プロキシ・キャッシュ 115 に格納されているデータに関する更新や変更を事前に確認することができる。オブジェクト（ディレクトリやファイルなど）の属性を NFS で得るためには、「GETATTR」コマンドが使用される（もちろん、他の適当な処理を用いてもよい）。各オブジェクトについて「GETATTR」コマンドを使用し、オブジェクトに変更があったか否かを示す属性をチェックする。一実施形態において、NFS トラフィックはメタデータの事前確認を行い、トラフィックを加速することができる。

【0097】

事前確認アクセスパターンには、例えば、シーケンシャル・アクセス、ランダム・アクセス、履歴を利用した隣接オブジェクトのリンクのアクセスなど、様々なパターンが使用される。シーケンシャル・アクセスによる確認とは、複数のオブジェクトの属性を順番に確認することを意味する。クライアント事前確認は、履歴によるリンクを作成し、クライアントは、構築された履歴によるリンクから利点を得る。履歴を利用した隣接オブジェクトのリンクについては、後で説明する。クライアントによっては、作成されたリンクに重みを加えてもよい。また、リンクを破壊する選択枝を追加することもできる。

【0098】

図 10 は、本発明の一実施形態による、事前確認アクセスを実施するためのプロキシ・キャッシュ 115 の構成要素を示すブロック図である。オブジェクト E0 は不可視のファイルハンドルを有する場合があるので、オブジェクト E0 の親ディレクトリを見つける方法は無い場合がある。プロキシ・キャッシュ 115 のメモリに格納される要素ハッシュテーブル 1005 およびディレクトリ・ハッシュテーブル 1010 は、ディレクトリ D に関する情報を有する。最初の 4 つの要素（例えば、E0 ~ E3）は、要素ハッシュテーブル 1010 にロードされる。この方法は、ディレクトリ D へのリンクを作成することができる、シーケンス E0 ~ E3 のキャプチャを短時間で行なうことができる。最初の 4 つのオブジェクト E0 ~ E3 は、オブジェクト E0 ファイルのファイルハンドルに基づいてテーブル 1010 内にハッシュされ、所定の隣接オブジェクト E1 ~ E3 および E0 は、要素ハッシュテーブル 1005 内のディレクトリ D を指すことになる（1015）。クライアント 105 はオブジェクト E0 の代わりにオブジェクト E1 または E2 を要求する場合があるので、オブジェクト E0 ~ E4 がロードされる。オブジェクト E0 エントリは、ハッシュテーブル 1005 にある親ディレクトリ D に関する十分な情報を有する。オブジェクト E0 がロードされると、次のポインタがマークされ、次のオブジェクトが E1 になる。このポインタは、アクセスパターンに関するシーケンシャル・リンクを作成する。図 1 のリンク方法を使用すると、ディレクトリ・クッキー順位の判定が不要になる。従って、このリンク方法は、クライアント 105 ごとに処理順序を管理することができる。テーブル 1010 にロードされるファイルは、「隠し」ファイルではない方が望ましい。隠しファイルとは、それを見るために特殊な方法を知っていなければならないファイルである。

【0099】

図 11 は、本発明の一実施形態による、確認後のオブジェクトを履歴的にリンクする方法を示すブロック図である。オブジェクトは、隣接オブジェクトのクライアントによる確認が済んでいれば、隣接オブジェクトの集合と対にされる（リンクされる）。例えば、クライアント 105 からの確認要求「GETATTR D」が、オブジェクト D を確認する。オブジェクト E0 および E1 は隣接オブジェクトであり、オブジェクト E0 および E1 は事前に確認されているので、オブジェクト D にリンクされる。

【0100】

オブジェクトの事前確認が特に有用である理由は、NFS プロトコルにおける「GETATTR」トラフィックが通常、ネットワークのボトルネックとなるからである。オブジェクトを事前確認することにより、WAN を介した NFS トラフィックを最適化することができる。これに対し、従来の方法で WAN を介した NFS トラフィックを最適化するためには、NFS トラフィックを適当なプロトコルに変換しなければならない。

【0101】

10

20

30

40

50

図12は、オブジェクトの属性の変化によってリンクが破壊されたところを示すブロック図である。例えば、クライアント105が「GETATTR E0」コマンドを送信した後、オブジェクトE0が事前確認され、オブジェクトE0の属性の変化が検出された場合、プロキシ・キャッシュ115は、サーバ110からオブジェクトE0を読み出し、オブジェクトE0をキャッシュすることにより処理を進める。

【0102】

プロトコル変換（変換要求）

図13は、クライアント105からのCIFSトラフィックに回答するキャッシュ・ヒット状態を説明するための、クライアント105、プロキシ・キャッシュ115、および、サーバ110を含むネットワークシステムを示す図である。クライアント105とプロキシ・キャッシュ115の間で使用されるプロトコルはCIFSプロトコルのようなオープン・プロトコルであり、プロキシ・キャッシュ115とサーバ110の間で使用されるプロトコルはNFSプロトコルであるものと仮定する。

10

【0103】

一実施形態においてプロキシ・キャッシュ115は、プロトコルレベルの変換を実施することなく、ファイルシステムレベルで、クライアント105からのCIFSトラフィックをサーバ110に対するNFSトラフィックに変換することができる（さらに、サーバ110からのNFSトラフィックをクライアント105に対するCIFSトラフィックに変換することもできる）。従って、プロキシ・キャッシュ115を用いると、CIFSクライアントは、キャッシュ・ミスがあっても、NFSサーバのデータにアクセスすることができる。一般に、CIFSプロトコルは、逐次処理であり、CIFSプロトコルに比べて処理回数が多いので、WANを介して使用することは実際的でない。さらに、CIFSプロトコルはあまり十分に定義されていないので、CIFSプロトコルをWANを介して使用するように最適化することは、プロトコル違反となる場合もある。これに対し、プロキシ・キャッシュ115は、クライアント105からのCIFSトラフィックをNFSトラフィックに変換することができるので、キャッシュ・ミスがあっても、サーバ110に対するアクセスが可能になる。

20

【0104】

図13は、クライアント105からのCIFS要求1305に関するキャッシュ・ヒット状態を示す図である。スタック1310は、要求105に回答してファイルシステム1315（例えば、Write-Anywhereファイルレイアウト、すなわちWAFLなど）に質問を送り、要求されたファイルをファイルシステム1315が有しているか否かを確認する。すると、ファイルシステム1315は、要求されたデータ1320をスタック1310に渡す。次に、スタック1310は、要求されたデータ1320を要求元クライアント105に送信する。同様の処理は、プロキシ・キャッシュ115がクライアントからNFS要求を受信し、キャッシュ・ヒット状態が発生した場合にも行なわれる。

30

【0105】

図14は、プロキシ・キャッシュ115においてキャッシュ・ミス状態を引き起こすCIFS要求の処理を示すブロック図である。スタック1310は、要求に応じて、要求されたデータに関する質問をファイルシステム1315に送る。ファイルシステム1315は、要求されたデータがファイルシステム1315内に存在しないので、NFSクライアントモジュール1325（このモジュールは通常、プログラムとして実施され、クライアント上で実行される）に質問を送る。すると、NFSクライアントモジュール1325はNFS要求1330をWANを介してサーバ110に送信する。サーバ110は、プロキシ・キャッシュ115からのNFS要求1330に回答して、要求されたデータをNFSプロトコルによりWANを介してNFSクライアントモジュール1325に送信する（1335）。次に、NFSクライアントモジュール1325は要求されたデータをファイルシステム1315に送信し、ファイルシステム1315はその要求されたデータをスタック1310に渡す。次にスタック1310は、要求されたデータ1340を要求元クライアント105に送信する。

40

50

【0106】

図15は、プロキシ・キャッシュ115にキャッシュ・ミス状態を引き起こすNFS要求1500の処理を示すブロック図である。スタック1510は、要求されたデータがファイルシステム1515内にあるか否かをチェックする。プロキシ・キャッシュ115は、キャッシュ・ミスに回答して、NFS要求1530をネットワーク・インタフェース1505からサーバ110に送信する。フィルエンジン1520は、サーバ応答1535を構文解析し、クライアント105とサーバ110の間の通信を可能にするデータ構造、および、不可視の特徴をもつファイルハンドルを使用する。要求されたデータ1525は、プロキシ・キャッシュ115に割り当てられた後、プロキシ・キャッシュ115からクライアント105に送信される。

10

【0107】

遅延書込み

図16は、遅延書込み処理を説明するための、本発明の一実施形態によるネットワークを示すブロック図である。プロキシ・キャッシュ115は、クライアント105からの書込み処理1605に回答し、遅延スパーズ書込み処理1610をサーバ110に対して実施するように構成される。例えば、クライアント105がプロキシ・キャッシュ115に対して書込み処理1605を実施する場合、プロキシ・キャッシュ115は、そのデータをサーバ110に書込み、サーバ110内のファイルを変更することができる。書込み処理1610は、プロキシ・キャッシュ115のバックグラウンド処理として実施してもよいし、後で実施してもよい。従って、書込みデータは最初にプロキシ・キャッシュ内のローカルディスクに送信され、リモートサーバに対する書込み処理は遅れる場合がある。

20

【0108】

プロキシ・キャッシュ115は、書込み処理1610を実施する前の準備領域としての働きを持ち、書込みデータをバッファリングしておくことができる。書込み処理1610は、書込み処理1605の速度よりも遅い速度で実施することができる。代替または追加として、書込み処理1610は、書込み処理1605よりも遅い時間に、すなわち書込み処理1605よりも後の時刻に実施してもよい。従って、プロキシ・キャッシュ115とサーバ110とを接続しているWANが故障した場合や、プロキシ・キャッシュ115とサーバ110との間に計画的切断が発生した場合、プロキシ・キャッシュ115は、WANが復旧した時点でサーバ110に対する書込み処理1610を実施することができる。この遅延書込み機能によれば、クライアント105の動作を低下させたり、クライアント105の動作に悪影響を及ぼすことなく、最終的には書込みデータをサーバ110に確実に格納することができる。従って、本発明の一実施形態によれば、サーバ110に対する書込み処理1610は通常自動的に実施されるので、データをコピーするための従来のソフトウェアは不要になる。

30

【0109】

書込み処理1605は、プロキシ・キャッシュ115に格納された読出し専用以外のデータを無効することができる。さらに、プロキシ・キャッシュ115にキャッシュされた書込みデータは、消去処理を必要としない。なぜなら、書込みデータはその後、例えば上記のようなFIFO処理を用いて、プロキシ・キャッシュ115から消去されるからである。

40

【0110】

プロキシ・キャッシュ115からサーバ110への遅延書込み処理は、次のように実施される。(1)即座に、(2)スケジュールに従って、(3)ネットワークのアイドル時間に、および/または、(4)帯域幅制限に応じて(例えば、書込み処理は、最大ネットワーク帯域幅が得られる時間に実施される)。プロキシ・キャッシュ115に対するトラフィックの帯域幅を制限するために、予約速度制御を行ってもよい。具体的には、帯域幅は、様々な要素に基づいて優先順位付けをすることによって制限または管理することができる。そうした要素には、例えば、(1)クライアントIP情報(レイヤ4 TCP/IP層情報)、(2)NFS情報(レイヤ7)(例えば、ユーザID、読出し/書込み等

50

の処理タイプなど)、(3)特定ファイル(例えば、特定のディレクトリの優先順位を他のディレクトリの優先順位よりも高くする場合がある)などがある。従って、帯域幅は、要求を送信しているクライアントのタイプや、上記のようなその他の要素に基づいて優先順位を付けることができる。

【0111】

Delegation(委譲)を用いた方法(Delegationは、サーバ内のオブジェクトを変更するパーミッションすなわち権利だけをプロキシ・キャッシュに与える)については、後で図18のところで説明する。Delegation技術は、CIFSおよびNFSバージョン4でサポートされている。

【0112】

図17は、サーバ110内のファイルシステム1705のうちの下位部分1703を含むアクティブデータセットをキャッシュするように構成されたプロキシ・キャッシュ115の一実施形態を示す。下位部分は、例えば「ルート」1707および「ユーザ1」の下にあるサブツリーである。サブツリー1703は、プロキシ・キャッシュ115とサーバ110の間の計画的切断に先立って、サーバ110からプロキシ・キャッシュ115にフェッチされる。従って、クライアント105は、プロキシ・キャッシュ115にアクセスすることにより、必須データまたは特定データを常に利用することができる。例えば、プロキシ・キャッシュ115は、計画的切断処理に先立って、サーバ110のサブツリー1706をフェッチするように構成される。従って、プロキシ・キャッシュは、従来の方法のミラーリング技術を使用せずにデータを複製する方法を提供することができるという利点を有する。

【0113】

プロキシ・キャッシュ115内にもともとキャッシュされていたデータのコピーは、サーバ110から全てのデータを取得するまで、サーバ110からの新たな更新データで置換されることはない。プロキシ・キャッシュ115におけるデータの置換は、ディレクトリレベルで実施してもよいし、ファイルレベルで実施してもよい。

【0114】

上で述べた名前空間の仮想化によれば、サーバ110内のオブジェクトの事前確認が可能になる。

【0115】

Delegation

図18は、本発明の一実施形態によるDelegationを使用するシステムを示すブロック図である。サーバ110(1)は、ディレクトリまたはファイル(複数の場合もあり)のコピー1810を有し、プロキシ・キャッシュ115は、論理コピー1810におけるデータのキャッシュされたコピー1805を有するものと仮定する。この例においてキャッシュされたコピー1805は、プロキシ・キャッシュ115内のアクティブデータセットに含まれることが好ましい。サーバ105は、論理コピー1810をプロキシ・キャッシュ115に「ロック」することができ、それによってプロキシ・キャッシュ115は、キャッシュ・ヒット時に、サーバ115の論理コピー1810にファイルの更新があったか否かをチェックすることなく、役目を果たすことができるようになる。サーバ110(1)は、サーバ110(1)の論理コピー1810に関するファイル更新をいずれも、ファイル更新ロックの所有者に通知する。プロキシ・キャッシュ115がサーバ内のファイルに対するロックを有している場合、そのファイルに対する書込み要求は通常いずれもプロキシ・キャッシュ115を通過し、その後プロキシ・キャッシュにより、そのファイルに対する遅延書込み処理が実施される。その後、プロキシ・キャッシュ115が所有するロックは、サーバによって無効にされる。

【0116】

クラスタリング

図19は、本発明の一実施形態によるクラスタ構成の複数のプロキシ・キャッシュ115a~115cを含むブロック図である。プロキシ・キャッシュの数は、図19のもの

10

20

30

40

50

は違っていてもよい。プロキシ・キャッシュのクラスタを用いると、一次プロキシ・キャッシュが故障したときに、クライアント要求トラフィックおよびデータトラフィックに関する代替経路の方法が可能になる。これらのプロキシ・キャッシュ115は、例えば、世界または地域の様々な部分に配置することができる。例えば、プロキシ・キャッシュ115aは一次キャッシュであり、その後これが故障した場合、他のプロキシ・キャッシュを一次プロキシ・キャッシュとして選択することにより、ネットワーク・トラフィックの代替経路を形成することができる。例えば、プロキシ・キャッシュ115bおよび115cが、一次プロキシ・キャッシュとして選択される。新たなプロキシ・キャッシュを選択する際の基準には例えば、サイズ(大きなサイズのプロキシ・キャッシュはキャッシュ・ヒット率を向上させる)、低減されるトラフィック待ち時間、位置(例えば、1つの「ピア」プロキシ・キャッシュにするか、近接する複数のプロキシ・キャッシュにするかなど)、および/または、要求されているデータなどがある。クライアント105は、グローバル名前空間機能およびファイルハンドル変換機能により、他のプロキシ・キャッシュの代わりに、1つのプロキシ・キャッシュにアクセスすることができる。種々のプロキシ・キャッシュはお互いを知らないので、ファイルハンドルの書換えを利用する。クライアントは、ファイルハンドルの書換えにより、自分の要求を特定のプロキシ・キャッシュに適切に宛先変更することができる。そして種々のプロキシ・キャッシュがそのメッセージを解釈する。これに対し、HTTPプロトコルでは、オブジェクトの名前がサーバの名前に既に埋め込まれている。

10

【0117】

20

図20は、本発明の一実施形態による、リバース・プロキシ構成に構成されたクライアント装置105(1)~105(n)、プロキシ・キャッシュ2205(1)~2205(p)およびサーバ2210を含むネットワーク200を示すブロック図である。プロキシ・キャッシュ(図22では、まとめて2205と呼ぶ)を用いると、サーバ2210に対する双方向トラフィックのスケーリングまたは分散が可能になる。プロキシ・キャッシュ2205は、LAN2215を介してサーバ2210にローカル接続される。変数pは任意の適当な整数値である。つまり、プロキシ・キャッシュ装置(図22では、まとめて2205と呼ぶ)の数は変更してもよい。例えば、サーバ2210のトラフィック処理速度が遅い場合、ネットワーク2200は、プロキシ・キャッシュ2205だけで実施することもできる。クライアント装置105は、WAN125を介してロードバランス調節スイッチのようなロード・バランサ2220に接続される。ロードバランス調節スイッチは、トラフィックをクライアント105とプロキシ・キャッシュ2205に分散させる。

30

【0118】

図22のリバース・プロキシ構成によれば、データをプロキシ・キャッシュ2205にオン・デマンドでキャッシュすることが可能になる。このオン・デマンドのキャッシュ処理は、標準的な複製技術に比べて効率が高く、標準複製技術の場合のように特殊なソフトウェアを必要としない。複数のプロキシ・キャッシュ2205を使用することで、サーバ2210に宛てられた要求を分散させることが可能になる。複数のプロキシ・キャッシュ2205を用いると、クライアント105のうちの1つからの要求に回答するキャッシュ・ヒットの可能性も全体的に上げることができる。なぜなら、各プロキシ・キャッシュ2205のアクティブデータセットが、要求されたデータを含む場合があるからである。従って、複数のプロキシ・キャッシュ2250を用いると、クライアント105からサーバ2210に送信される要求の量を減らすことができる。

40

【0119】

利点として、サーバ2210の処理速度が遅い場合は、複数のプロキシ・キャッシュ2205を使用することで、キャッシュ・ヒットの可能性を向上させることができ、クライアント105に対してオン・デマンドのデータキャッシュ機能を提供することができる。従って、サーバ2210の処理速度が遅いことが、クライアント要求のためのボトルネックにはなくなる。図22のリバース・プロキシ構成は、読出しトラフィックに(例えば、動画のグラフィック・レンダリングなどの用途に)関して非常に有用である。また、

50

リバース・プロキシ構成の場合、ファイルハンドルの変形は一般に必要とされない。さらに、プロキシ・キャッシュの使用はファイルサーバの使用に比べて安価であり、また、プロキシ・キャッシュにおけるデータの複製は特殊なソフトウェアを必要とせずに自動的に実施することができる。

【0120】

ロード・バランサ2220は、キャッシュ・ヒットからの要求およびデータを分散させる。キャッシュ・ミスがあった場合（すなわち、要求されたデータをアクティブデータセット内に有しているプロキシ・キャッシュ205が無い場合）、プロキシ・キャッシュ205のうちの1つは、要求されたデータを求める要求をサーバ210に送り、そのデータのキャッシュされたコピーを要求元サーバに返送する。

10

【0121】

図21は、単一のプロキシ・キャッシュ2205(1)がサーバ2210にリバース・プロキシ構成で接続されたネットワーク2300を示すブロック図である。この例では、サーバ2210の処理速度が十分に遅いものになっている点に注意して欲しい。従って、ネットワーク2300で実施される単一のプロキシ・キャッシュ2205は、クライアント要求を受信し、キャッシュ・ヒットがあった場合はデータを要求元クライアント105に返送し、キャッシュ・ミスがあった場合はサーバ2210にアクセスする。

【0122】

図22は、本発明の他の実施形態によるネットワーク2400を示すブロック図である。プロキシ・キャッシュ2205(1)~2205(p)は、例えばネットワーク・スイッチ2405を介して、複数のサーバ2210(1)~2210(m)と通信することができる。変数mは、任意の適当な整数値である。通常、ネットワーク・スイッチ2405はLAN2215内で実施される。従って、プロキシ・キャッシュ2205は、複数のサーバ2210に対する双方向のトラフィックをスケーリングまたは分散させることができ、サーバ2210に対する双方向のトラフィックを加速させることができる。図22のリバース・プロキシ構成は、特にサーバの処理能力が低い場合に、サーバ2210がクライアント・トラフィックに対するボトルネックとして作用することを防止する。

20

【0123】

図23は、本発明のさらに他の実施形態によるネットワーク2450を示すブロック図である。プロキシ・キャッシュ2505は、プロキシ・キャッシュ2205(1)~2250(p)と一緒にクラスタ化され、クラスタ化プロキシ・キャッシュについて上に記載した種々の機能を実施することができる。

30

【0124】

本発明の実施形態は、全ての要求についてクライアントがサーバにアクセスしなければならない従来の方法とは違い、クライアントの待ち時間を低減することができるという利点を有する。本発明の実施形態は、クライアントに特殊なソフトウェアを必要とせず、遠くの場合からサーバを簡単に管理することができるので、クライアントおよび/またはサーバの管理の容易さを向上させるという利点を有する。本発明の実施形態は、従来の方法のミラーリングステップを不要にすることにより、データ複製の容易さを向上させることができるという利点を有する。具体的には、本発明の実施形態は、所与の時間間隔でデータのミラーリングを行なうことが可能な管理者を必要としない。本発明の実施形態は、アクティブデータセット全体のサイズがサーバ内のデータセット全体のサイズよりも小さいので、必要とされるディスク要件も少なく済むという利点を有する。ディスク要件が少ないことから、コストも少なく済み、故障問題も少なくなる。本発明の実施形態は、多数のクライアントがリバース・プロキシ構成のネットワークを飽和させることなく、多数のデータ要求を行なえるという利点を有する。

40

【0125】

また、本明細書に記載する種々のエンジンまたはモジュールは、例えば、ソフトウェア、コマンド、データファイル、プログラム、コード、モジュール、命令などであってもよく、適当なメカニズムをさらに含むものであってもよい。

50

【0126】

本明細書全体を通じて「1つの実施形態」、「一実施形態」または「特定の実施形態」と呼ばれるものは、本発明の少なくとも1つの実施形態に含まれる実施形態に関連して説明された特定の特徴、構造または特性を意味している。従って、本明細書全体を通じて種々の場所に記載されている「1つの実施形態において」、「一実施形態において」または「特定の実施形態において」のような語句は、すべて同じ実施形態を指しているとは必ずしも限らない。さらに、特定の特徴、構造または特性は、任意の適当な方法で1以上の実施形態と組み合わせることもできる。

【0127】

上に記載された教示を参照すれば、上記の実施形態および方法に対する変形および変更は他にも可能である。また、本発明の一実施形態に関する少なくとも一部の構成要素は、相互接続された構成要素および回路のネットワークを用いて、特定用途向け集積回路、プログラマブル・ロジック・デバイスまたはフィールド・プログラマブル・ゲート・アレイを用いて、プログラムされた汎用デジタルコンピュータを用いて実施される。接続は、有線でも、無線でも、モデムによるものなどであってもよい。

【0128】

また、図面/図に描かれた要素のうち1以上は、特定の用途にとって有用となるように、さらに離れた態様で実施することも、さらに一体化した態様で実施することもでき、場合によっては、除去したり、動作しないようにすることもできる。

【0129】

さらに、機械読取可能媒体に格納可能なプログラムまたはコードを実施し、上記の方法のうちのいずれかをコンピュータで実施することも、本発明の範囲内である。

【0130】

さらに、図/図面中の信号矢印は、例として解釈すべきものであり、特にことわりがない限り、限定の意味で解釈してはならない。さらに、本明細書の開示に使用される「または」という用語は通常、特にことわりがない限り、「および/または」の意味で使用している。また、構成要素またはステップの組み合わせも考えられる。その場合、分離または結合する能力の解釈が不明瞭であるので、専門用語は予測される。

【0131】

本明細書の説明および特許請求の範囲の記載に使用されている「a」、「an」および「the」は、文脈から明らかにそうでないと言えない限り、複数の場合も含む。また、本明細書の説明および特許請求の範囲の記載に使用されている「in」は、文脈から巻きからかにそうでないと言えない限り、「in」の意味と「on」の意味とを含む。

【0132】

さらに本発明の態様は、ハードウェア、ソフトウェア、ファームウェア、および/または、それらの組み合わせで実施することができる。

【0133】

要約書に記載したものを含む本発明の例示的な実施形態に関する上記の説明は、本発明を網羅的に記載しようとするものでもなければ、本発明を開示した形態に厳密に制限しようとするものでもない。本明細書には、例示の目的で本発明の具体的実施形態や例が記載されているが、本発明の範囲内で種々の均等な変更が可能であることは、当業者にとって明らかであろう。

【0134】

上記の詳細な説明を参照すれば、本発明に対してそれらの変更を行なうことができる。特許請求の範囲に使用されている用語は、本発明を明細書および特許請求の範囲に開示されている特定の実施形態に制限するものとして解釈してはならない。そうではなく、本発明の範囲は特許請求の範囲によって完全に決定され、特許請求の範囲は、特許請求の範囲に関する確立された解釈の原則に従って解釈しなければならない。

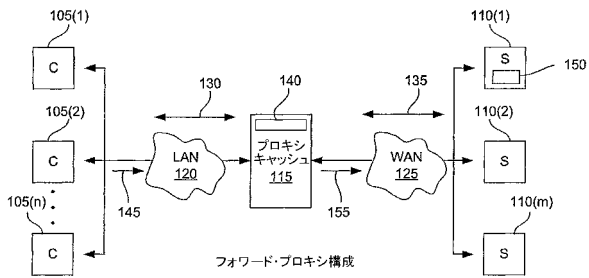
【図面の簡単な説明】

【0135】

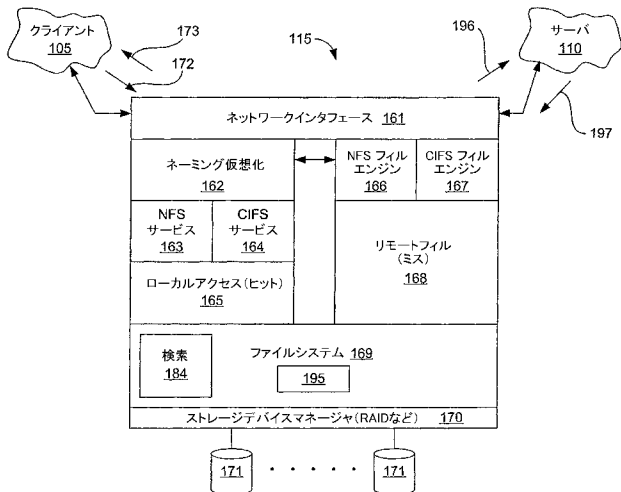
- 【図 1 A】本発明の一実施形態による装置を示すブロック図である。
- 【図 1 B】本発明の一実施形態によるプロキシ・キャッシュを示すブロック図である。
- 【図 1 C】キャッシュ・ヒット処理およびキャッシュ・ミス処理を示すブロック図である。
- 【図 2】本発明の一実施形態による加齢方式または重要度方式を用いてアクティブデータセットを管理する方法を示すブロック図である。
- 【図 3】本発明の実施形態によりアクティブデータセットを管理する他の方法を示すブロック図である。
- 【図 4】ファイルハンドルを示すブロック図である。
- 【図 5】各サーバ内のボリュームの例を示すブロック図である。 10
- 【図 6】本発明の一実施形態による、特定サーバにおける各ボリュームの F S i d 値のマッピングを示すブロック図である。
- 【図 7】本発明の一実施形態による、同一にマッピングされたファイルハンドルを有する複数のプロキシ・キャッシュを備えたネットワークシステムを示すブロック図である。
- 【図 8 A】本発明の一実施形態による、仮想名前空間を可能にするためのローカルディレクトリを示すブロック図である。
- 【図 8 B】1つのマウントポイントに一連の属性が関連付けられた複数のマウントポイントを作成する方法を示すブロック図である。
- 【図 8 C】キャッシュデータの整合性をとる方法を示すブロック図である。
- 【図 8 D】キャッシュデータの整合性をとる方法を示すブロック図である。 20
- 【図 8 E】キャッシュデータの整合性をとる方法を示すブロック図である。
- 【図 9】本発明の一実施形態による、サーバ内のデータの属性のプリフェッチまたは事前確認を示すブロック図である。
- 【図 10】本発明の一実施形態による事前確認アクセスを実施するための、プロキシ・キャッシュ内の構成要素を示すブロック図である。
- 【図 11】本発明の一実施形態による、確認されたオブジェクトを履歴によりリンクさせる方法を示すブロック図である。
- 【図 12】本発明の一実施形態による、オブジェクトの属性の変化に応じてリンクが破壊されたところを示すブロック図である。
- 【図 13】本発明の一実施形態による、クライアントからの C I F S トラフィックにตอบสนองするキャッシュヒット状態を説明するための、クライアント、プロキシ・キャッシュおよびサーバを含むネットワークシステムを示す図である。 30
- 【図 14】本発明の一実施形態による、プロキシ・キャッシュにキャッシュ・ミス状態を引き起こす C I F S 要求の処理を示すブロック図である。
- 【図 15】本発明の一実施形態による、プロキシ・キャッシュにキャッシュ・ミス状態を引き起こす N F S 要求の処理を示すブロック図である。
- 【図 16】遅延書込みを説明するための、本発明の一実施形態によるネットワークを示すブロック図である。
- 【図 17】サーバの所定の記憶領域にあるデータの特定の論理コピーのキャッシュされたコピーを含むアクティブデータセットをキャッシュするように構成されたプロキシ・キャッシュの一実施形態を示す図である。 40
- 【図 18】本発明の一実施形態によるデータコピーレンシを提供するシステムのブロック図である。
- 【図 19】本発明の一実施形態による、クラスタ構成の複数のプロキシ・キャッシュを含むブロック図である。
- 【図 20】本発明の一実施形態による、リバース・プロキシ構成のクライアント、プロキシ・キャッシュおよびサーバを含むネットワークを示すブロック図である。
- 【図 21】単一のプロキシキャッシュ 205 (1) がサーバ 210 にリバース・プロキシ構成で接続されたネットワーク 300 を示すブロック図である。
- 【図 22】本発明の他の実施形態によるネットワーク 2400 を示すブロック図である。 50

【図23】本発明の他の実施形態によるネットワーク2500を示すブロック図である。

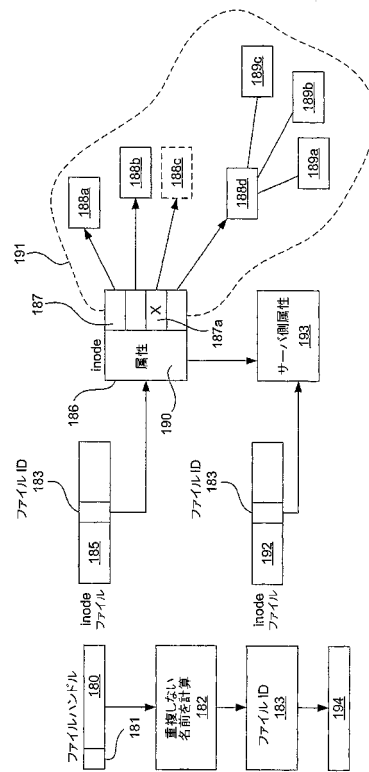
【図1A】



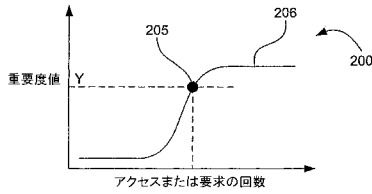
【図1B】



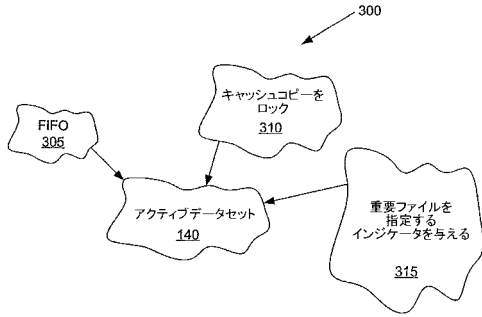
【図1C】



【 図 2 】



【 図 3 】

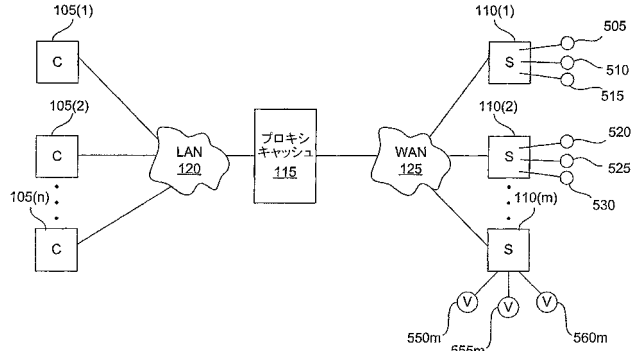


【 図 4 】

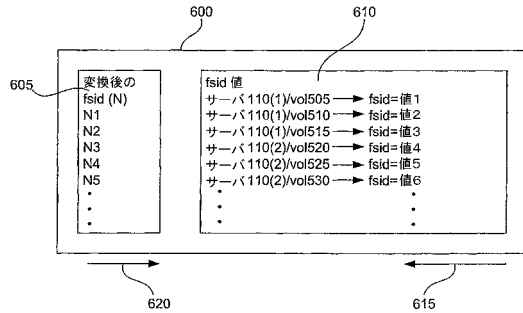
ファイルハンドル

fsid	ファイル ID	生成番号	その他のデータ
405	410	415	420

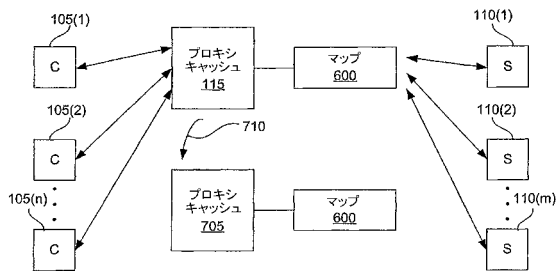
【 図 5 】



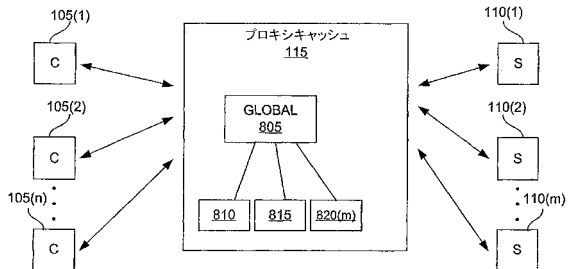
【 図 6 】



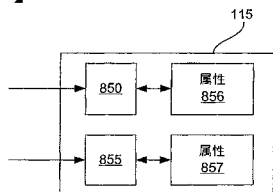
【 図 7 】



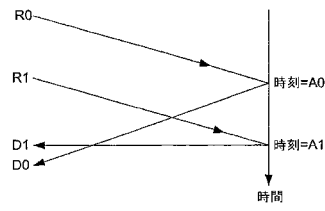
【 図 8 A 】



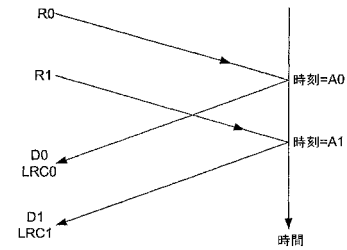
【 図 8 B 】



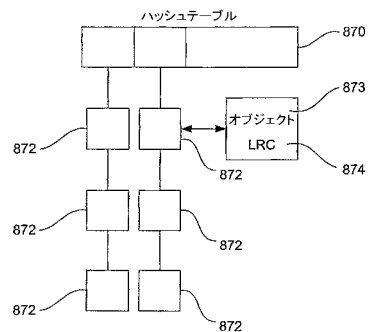
【 図 8 C 】



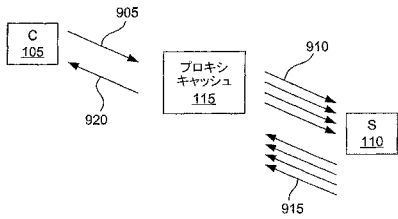
【 図 8 D 】



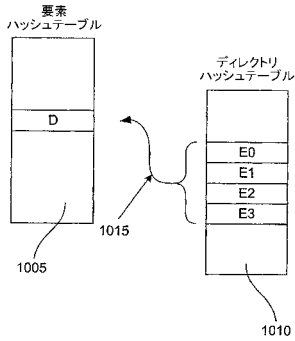
【 図 8 E 】



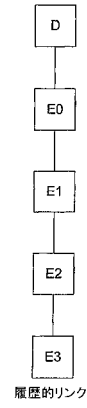
【図9】



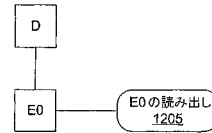
【図10】



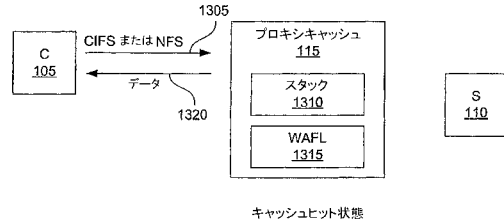
【図11】



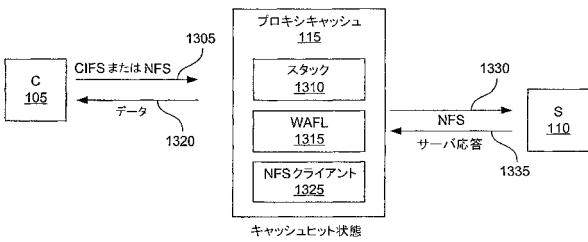
【図12】



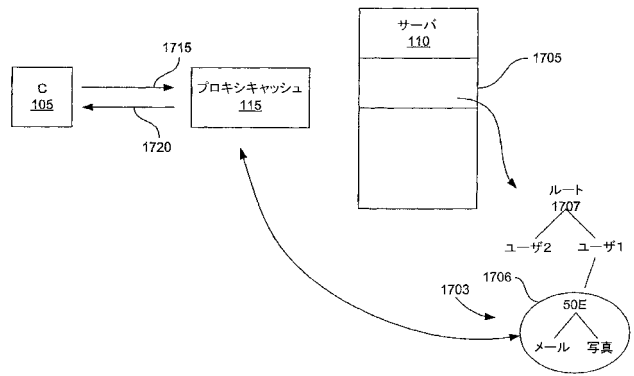
【図13】



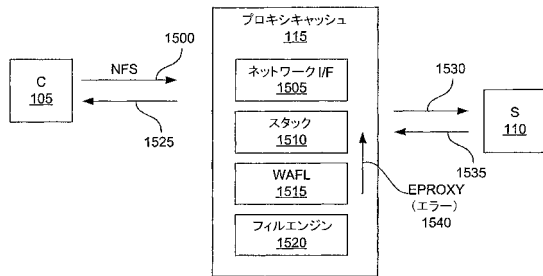
【図14】



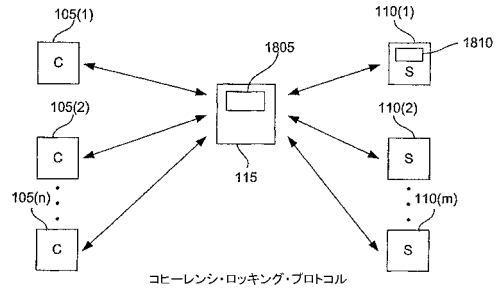
【図17】



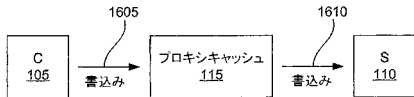
【図15】



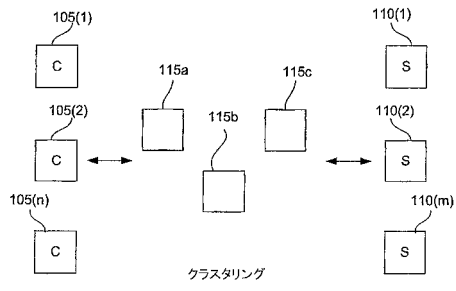
【図18】



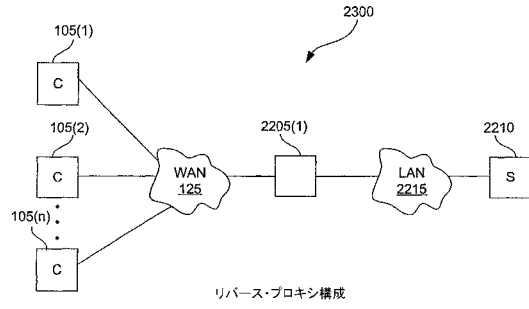
【図16】



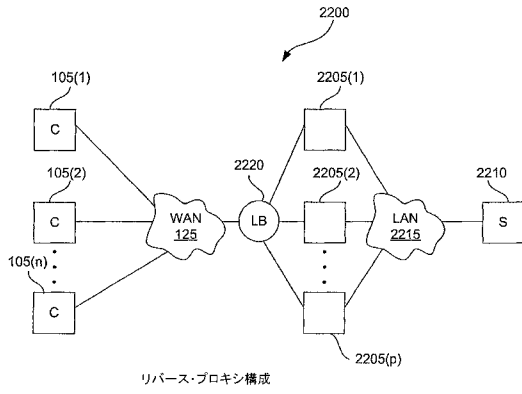
【図19】



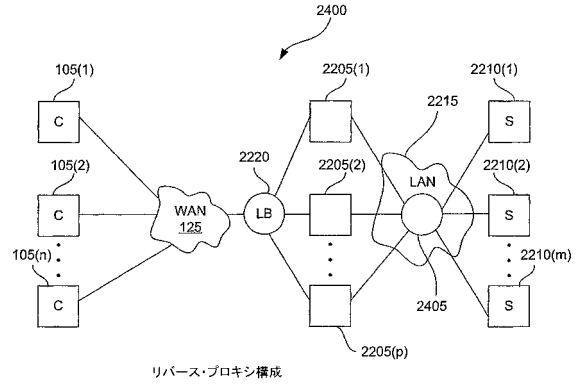
【図21】



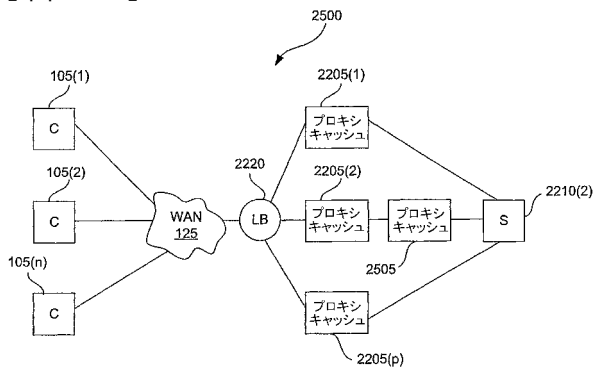
【図20】



【図22】



【図23】



【手続補正書】

【提出日】平成17年7月11日(2005.7.11)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

データをネットワーク内にキャッシュする装置であって、
クライアントからオブジェクトの要求を受信し、サーバからデータブロックをフェッチするように構成されたプロキシ・キャッシュを含み、

前記プロキシ・キャッシュは、階層構造の前記データブロックをアクティブデータセットとして前記オブジェクト内にキャッシュするように構成される、装置。

【請求項2】

前記データブロックはバッファ・ツリーに入れられ、ファイル内に前記階層関係が可能である、請求項1に記載の装置。

【請求項3】

クライアントからの要求に回答して、前記アクティブデータセット中にデータブロックが配置される、請求項1に記載の装置。

【請求項4】

前記アクティブデータセットは、少なくとも前記サーバに格納されたデータの一部を含む、請求項1に記載の装置。

【請求項5】

前記アクティブデータセットは、特定の時間によって定義される、請求項1に記載の装置。

【請求項6】

前記プロキシ・キャッシュは、前記サーバによって要求されたファイルのデータブロックが欠けていた場合、前記クライアントからの要求に回答して該データブロックを前記サーバからフェッチするように構成される、請求項1に記載の装置。

【請求項7】

前記プロキシ・キャッシュは、前記サーバによって要求されたファイルが前記プロキシ・キャッシュに存在しない場合、前記クライアントから前記サーバへの要求に回答して前記サーバからデータブロックをフェッチするように構成される、請求項1に記載の装置。

【請求項8】

前記プロキシ・キャッシュは、クライアントによって要求されたデータファイル中のデータブロックの有無を判定する計算の一部として重複しない名前を計算する、請求項1に記載の装置。

【請求項9】

前記オブジェクトはデータファイルを含む、請求項1に記載の装置。

【請求項10】

前記オブジェクトはファイルの一部を含む、請求項1に記載の装置。

【請求項11】

前記オブジェクトはディレクトリを含む、請求項1に記載の装置。

【請求項12】

データをネットワーク内にキャッシュする方法であって、
クライアントからプロキシ・キャッシュへのオブジェクトの要求を受信するステップと、

要求された前記オブジェクトが前記プロキシ・キャッシュ内にある場合、要求された前記オブジェクトを前記クライアントに送信するステップと、

要求された前記オブジェクト中のデータブロックが欠けていた場合、該データブロックをサーバからフェッチし、該データブロックをアクティブデータセットとして階層関係に入れ、要求された前記オブジェクトを前記クライアントに送信するステップと、
からなる方法。

【請求項 13】

前記サーバによって要求されたファイルが存在しない場合、前記クライアントから前記サーバへの要求に回答して前記サーバからデータブロックをフェッチするステップをさらに含む、請求項 12 に記載の方法。

【請求項 14】

前記データブロックがバッファ・ツリーに入れられ、ファイル内に前記階層関係が可能である、請求項 12 に記載の方法。

【請求項 15】

前記アクティブデータセットは、少なくとも前記サーバに格納されたデータの一部を含む、請求項 12 に記載の方法。

【請求項 16】

前記アクティブデータセットは、特定の時間によって定義される、請求項 12 に記載の方法。

【請求項 17】

クライアントによって要求されたデータファイル中のデータブロックの有無を判定する計算の一部として重複しない名前を計算するステップを含む、請求項 12 に記載の方法。

【請求項 18】

置換ポリシーに基づいて前記アクティブデータセットの一部を消去するステップをさらに含む、請求項 12 に記載の方法。

【請求項 19】

オブジェクトは、ファイルの少なくとも一部を含む、請求項 12 に記載の方法。

【請求項 20】

データをネットワーク内にキャッシュする方法であって、
クライアントからプロキシ・キャッシュへのオブジェクトの要求を受信するステップと、

要求された前記オブジェクトを前記クライアントに送信するステップと、

前記データブロックをアクティブデータセットの一部として階層関係に入れ、要求された前記オブジェクトを前記クライアントに送信するステップと、

からなる方法。

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/US03/29398		
A. CLASSIFICATION OF SUBJECT MATTER				
IPC(7) : G06F 15/167, 173 US CL : 709/225, 216				
According to International Patent Classification (IPC) or to both national classification and IPC				
B. FIELDS SEARCHED				
Minimum documentation searched (classification system followed by classification symbols) U.S. : 709/225, 216				
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched				
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EAST				
C. DOCUMENTS CONSIDERED TO BE RELEVANT				
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.		
Y	US 6,311,216 B1(SMITH et al.) 30 October 2001, col. 1, line 10 - col. 21, line 13.	1-63		
Y	US 6,438,652 B1 (JORDAN et al) 20 August 2002, col. 1, line 13 to col. 9, line 58.	1-63		
Y	US 6,442,651 B2 (CROW et al.) 27 August 2002, col. 1, line 13 to col. 5, line 55.	1-63		
A	US 6,178,461 B1 (CHAN et al) 23 January 2001, col. 1, line 15 to col. 11, line 63.	1-63		
A, P	US 6,553,411 B1 (DIAS et al) 22 April 2003, col. 1, line 13 to col. 6, line 15	1-63		
A	US 2002/0026560 A1 (JORDAN et al) 28 February 2002, [0002-0041].	1-63		
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.				
* Special categories of cited documents: <table border="0" style="width: 100%;"> <tr> <td style="width: 50%; vertical-align: top;"> <p>"A" document defining the general state of the art which is not considered to be of particular relevance</p> <p>"E" earlier application or patent published on or after the international filing date</p> <p>"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)</p> <p>"O" document referring to an oral disclosure, use, exhibition or other means</p> <p>"P" document published prior to the international filing date but later than the priority date claimed</p> </td> <td style="width: 50%; vertical-align: top;"> <p>"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</p> <p>"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone</p> <p>"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art</p> <p>"&" document member of the same patent family</p> </td> </tr> </table>			<p>"A" document defining the general state of the art which is not considered to be of particular relevance</p> <p>"E" earlier application or patent published on or after the international filing date</p> <p>"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)</p> <p>"O" document referring to an oral disclosure, use, exhibition or other means</p> <p>"P" document published prior to the international filing date but later than the priority date claimed</p>	<p>"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</p> <p>"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone</p> <p>"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art</p> <p>"&" document member of the same patent family</p>
<p>"A" document defining the general state of the art which is not considered to be of particular relevance</p> <p>"E" earlier application or patent published on or after the international filing date</p> <p>"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)</p> <p>"O" document referring to an oral disclosure, use, exhibition or other means</p> <p>"P" document published prior to the international filing date but later than the priority date claimed</p>	<p>"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</p> <p>"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone</p> <p>"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art</p> <p>"&" document member of the same patent family</p>			
Date of the actual completion of the international search 07 April 2004 (07.04.2004)		Date of mailing of the international search report 19 APR 2004		
Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US Commissioner for Patents P.O. Box 1450 Alexandria, Virginia 22313-1450 Facsimile No. (703) 305-3230		Authorized officer David A. Wiley Telephone No. 703-305-3900		

フロントページの続き

(81) 指定国 AP(GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VC, VN, YU, ZA, ZM, ZW

(72) 発明者 アカウィー, エマニュエル
アメリカ合衆国カリフォルニア州 9 4 1 0 7, サンフランシスコ, コネチカット・ストリート・5
7 5

(72) 発明者 プラカシュ, アシシュ
アメリカ合衆国ノースカロライナ州 2 7 5 6 0, モリスヴィル, ラ・ジョラ・レーン・1 0 6

(72) 発明者 アムダー, マシュー
アメリカ合衆国カリフォルニア州 9 4 1 0 7, サンフランシスコ, ミシシッピー・ストリート・5
8 0

(72) 発明者 アイヤー, カーティク
アメリカ合衆国カリフォルニア州 9 5 0 5 1, サンタクララ, ホメステッド・ロード・3 1 3 1,
アパートメント・ナンバー 2 3 エイチ

Fターム(参考) 5B082 FA12 HA02 HA05 HA08