

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2017-123040

(P2017-123040A)

(43) 公開日 平成29年7月13日(2017.7.13)

(51) Int. Cl.	F 1	テーマコード (参考)
<b>G06F 12/00 (2006.01)</b>	G06F 12/00 545F	
	G06F 12/00 514R	
	G06F 12/00 520J	

審査請求 未請求 請求項の数 10 ○L (全 33 頁)

(21) 出願番号 特願2016-1571 (P2016-1571)  
 (22) 出願日 平成28年1月7日(2016.1.7)

(71) 出願人 000004237  
 日本電気株式会社  
 東京都港区芝五丁目7番1号  
 (74) 代理人 100109313  
 弁理士 机 昌彦  
 (74) 代理人 100124154  
 弁理士 下坂 直樹  
 (72) 発明者 大谷 敦久  
 東京都港区芝五丁目7番1号  
 日本電気株式会社内

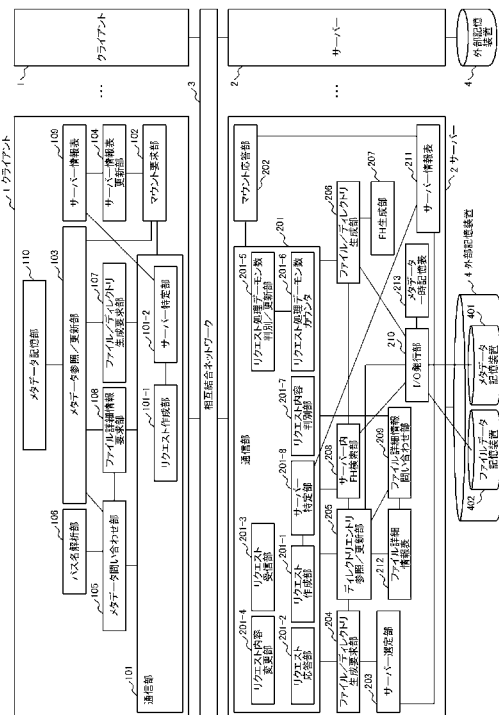
(54) 【発明の名称】 サーバー装置、分散ファイルシステム、分散ファイルシステム制御方法、および、プログラム

(57) 【要約】

【課題】 ファイルデータもメタデータも複数のサーバー装置に分散配置されているシステムにおいて、クライアント装置による効率的なアクセス対象サーバー装置の特定と、多数のリクエストが特定のサーバー装置に集中した際の効率的なアクセスを可能とする。

【解決手段】 ファイルが存在するディレクトリを管理するサーバー装置は、オープン対象ファイルのファイルハンドルと、ファイル詳細情報のリクエストを1つにまとめて受信し、同時に処理する。この際、当該サーバー装置は、ファイル詳細情報をメモリ上にキャッシュしておく。当該サーバー装置は、所定閾値以上のリクエストを受信した場合、リクエストを1つにまとめて処理することを止める。この場合、クライアント装置が、個別にファイル詳細情報のリクエストを、ファイルを管理するサーバー装置に送信する。

【選択図】 図2



## 【特許請求の範囲】

## 【請求項 1】

ファイルを記憶するファイルデータ記憶装置と、

ディレクトリ階層における直下のディレクトリ、または、ファイルの、a) 識別子であって、当該ディレクトリ、または、当該ファイルを管理するサーバー装置の識別情報であるサーバーIDを含むファイルハンドル(以降、FHと略記)、b) 名前、および、c) ファイルまたはディレクトリの区別を示すファイルタイプを関連付けて記憶するディレクトリ、および、前記ファイルデータ記憶装置に記憶されているファイルのファイル詳細情報を記憶するメタデータ記憶装置と、に接続され、

ファイル詳細情報のキャッシュ用メモリ領域であるファイル詳細情報表と、

クライアント装置、若しくは、他の前記サーバー装置から処理リクエストを受信、または、他の前記サーバー装置にリクエストを送信するとともに、処理中または送信中のリクエスト数をカウントする通信部と、

ファイルのFHを入力されると、FHが包含するサーバーIDの前記サーバー装置にファイル詳細情報を要求するリクエストを送信して、ファイル詳細情報を受信するファイル詳細情報問い合わせ部と、

クライアント装置から、ディレクトリ階層における直下のファイルの名前を含む、1) ファイルのFH、および、ファイルタイプにたいする要求、並びに、2) ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストが受信されたとき、a) 前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のファイルのFHとファイルタイプを取得し、b 1) 前記通信部がカウントしたリクエスト数が所定閾値を超えておらず、c 1) 受信された名前のファイルのファイル詳細情報が前記ファイル詳細情報表にキャッシュされていれば、キャッシュされているファイル詳細情報を取得し、c 2) 前記ファイル詳細情報表にキャッシュされていなければ、取得したファイルのFHを前記ファイル詳細情報問い合わせ部に入力してファイル詳細情報を取得して、前記ファイル詳細情報表にキャッシュし、取得したファイルのFH、ファイルタイプとファイル詳細情報とを前記クライアント装置に返信し、b 2) 前記通信部がカウントしたリクエスト数が所定閾値を超えている場合は、入力されたリクエストを変更してファイル詳細情報の出力を中止し、取得したファイルのFH、ファイルタイプを前記クライアント装置に返信し、また、前記サーバー装置から返信されたファイルのFHを、改めてファイルを管理する前記サーバー装置に送信する前記クライアント装置から、または、他の前記サーバー装置の前記ファイル詳細情報問い合わせ部から、前記ファイルデータ記憶装置に記憶されているファイル、のFHを含むリクエストが受信されると、前記メタデータ記憶装置から、入力されたFHを識別子とするファイルのファイル詳細情報を取得して返信するディレクトリエントリ参照/更新部と、を備えたサーバー装置。

## 【請求項 2】

前記サーバー装置は、さらに、

前記クライアント装置から、新たなディレクトリ、または、ファイルの第1の生成要求が送信されたとき、前記第1の生成要求に含まれる名前に依存しない方法で新たなディレクトリ、または、ファイルを管理する前記サーバー装置を選択するサーバー選定部と、

前記サーバー選定部が選択した前記サーバー装置に新たなディレクトリ、または、ファイルの第2の生成要求を送信して、生成された新たなディレクトリ、または、ファイルのFHを受信し、さらに、新たなファイルを生成した時はファイル詳細情報も受信する、ファイル/ディレクトリ生成要求部と、

前記第2の生成要求を受信すると、a) ディレクトリを生成して前記メタデータ記憶装置に格納、または、b) ファイル詳細情報を生成して前記メタデータ記憶装置に格納すると共に前記第2の生成要求送信元に返信し、さらに、自装置のサーバーID及び自装置内一意の値から新たなディレクトリ、または、ファイルのFHを作成して、作成したFHを前記メタデータ記憶装置に格納すると共に前記第2の生成要求送信元に返信するファイル/ディレクトリ生成部とを、備え、

10

20

30

40

50

前記ディレクトリエントリ参照/更新部は、ファイル/ディレクトリ生成要求部が受信したFH、並びに、前記第1の生成要求に含まれる名前、および、ファイルタイプを前記メタデータ記憶装置に格納されているディレクトリに記憶し、さらに、受信した場合はファイル詳細情報を前記ファイル詳細情報表に記憶する、請求項1のサーバー装置。

【請求項3】

前記ディレクトリエントリ参照/更新部は、前記クライアント装置から、ディレクトリ階層における直下のディレクトリの名前を入力されると、前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のFHとファイルタイプを取得して出力する、請求項1乃至請求項2の何れか1項のサーバー装置と、

まず、1)上位のディレクトリを管理する前記サーバー装置の前記ディレクトリエントリ参照/更新部に、直下のディレクトリの名前を入力して、直下のディレクトリのFHとファイルタイプを取得することを、前記パス名に名前が含まれるディレクトリについて、ルートディレクトリを管理する予め定められた前記サーバー装置を起点に、上位から順に繰り返すディレクトリ探索を実行して、パス名中の最下位ディレクトリのFHとファイルタイプを取得し、次に、2)最下位ディレクトリのFHが包含するサーバーIDの前記サーバー装置の前記ディレクトリエントリ参照/更新部に、2-1)ファイルのFH、および、ファイルタイプにたいする要求、並びに、2-2)ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストを送信して、a)前記パス名のファイルのFH、ファイルタイプ、および、ファイル詳細情報を取得する、または、b)前記パス名のファイルのFH、および、ファイルタイプを取得するメタデータ問い合わせ部と、

前記メタデータ問い合わせ部がファイル詳細情報を取得できなかったとき、前記メタデータ問い合わせ部が取得した前記パス名のファイルのFHを、改めてファイルを管理する前記サーバー装置に送信して、ファイル詳細情報を取得するファイル詳細情報要求部と、を備える前記クライアント装置と、を包含する分散ファイルシステム。

【請求項4】

前記クライアント装置は、

ディレクトリの、名前、FH、及び、ファイルタイプを関連付けてキャッシュするメタデータ記憶部を、さらに、備え、

前記メタデータ問い合わせ部は、前記ディレクトリ探索の過程で前記パス名から、ディレクトリの名前を取り出すと、a1)取り出した名前(以降、取得ディレクトリ名)が前記メタデータ記憶部にキャッシュされていなければ、前記取得ディレクトリ名のディレクトリの上位のディレクトリを管理する前記サーバー装置の前記ディレクトリエントリ参照/更新部に前記取得ディレクトリ名を送信して、前記取得ディレクトリ名のディレクトリのFHとファイルタイプを得るとともに、前記メタデータ記憶部に、前記取得ディレクトリ名、並びに、前記サーバー装置から取得したFHとファイルタイプをキャッシュし、a2)キャッシュされていれば前記メタデータ記憶部から前記取得ディレクトリ名のディレクトリのFHとファイルタイプを取得する、請求項2乃至請求項3の何れか1項の分散ファイルシステム。

【請求項5】

前記メタデータ記憶部は、ファイルの、名前、FH、ファイルタイプ、及び、ファイル詳細情報を関連付けて、さらに、キャッシュし、

前記メタデータ問い合わせ部は、前記パス名から、ファイルの名前を取り出すと、a1)取り出した名前(以降、取得ファイル名)が前記メタデータ記憶部にキャッシュされていなければ、前記サーバー装置から、前記取得ファイル名のファイルの、FH、ファイルタイプ、及び、ファイル詳細情報を得るとともに、前記メタデータ記憶部に、前記取得ファイル名、並びに、前記サーバー装置から取得したFH、ファイルタイプ、及び、ファイル詳細情報をキャッシュし、a2)キャッシュされていれば前記メタデータ記憶部から前記取得ファイル名のファイルのFH、ファイルタイプ、及び、ファイル詳細情報を取得する、請求項4の分散ファイルシステム。

【請求項6】

ファイルを記憶するファイルデータ記憶装置と、

ディレクトリ階層における直下のディレクトリ、または、ファイルの、a) 識別子であって、当該ディレクトリ、または、当該ファイルを管理するサーバー装置の識別情報であるサーバーIDを含むファイルハンドル(以降、FHと略記)、b) 名前、および、c) ファイルまたはディレクトリの区別を示すファイルタイプを関連付けて記憶するディレクトリ、および、前記ファイルデータ記憶装置に記憶されているファイルのファイル詳細情報を記憶するメタデータ記憶装置と、に接続され、ファイル詳細情報のキャッシュ用メモリ領域であるファイル詳細情報表を備える前記サーバー装置が、

クライアント装置、若しくは、他の前記サーバー装置から処理リクエストを受信、または、他の前記サーバー装置にリクエストを送信するとともに、処理中または送信中のリクエスト数をカウントし、

クライアント装置から、ディレクトリ階層における直下のファイルの名前を含む、1) ファイルのFH、および、ファイルタイプにたいする要求、並びに、2) ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストを受信すると、a) 前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のファイルのFHとファイルタイプを取得し、b1) カウントしたリクエスト数が所定閾値を超えておらず、c1) 受信された名前のファイルのファイル詳細情報が前記ファイル詳細情報表にキャッシュされていれば、キャッシュされているファイル詳細情報を取得し、c2) 前記ファイル詳細情報表にキャッシュされていなければ、取得したファイルのFHを含むリクエストを、当該FHが包含するサーバーIDの前記サーバー装置に送信してファイル詳細情報を取得して、前記ファイル詳細情報表にキャッシュし、取得したファイルのFH、ファイルタイプとファイル詳細情報とを前記クライアント装置に返信し、b2) カウントしたリクエスト数が所定閾値を超えている場合は、入力されたリクエストを変更してファイル詳細情報の出力を中止し、取得したファイルのFH、ファイルタイプを前記クライアント装置に返信し、

また、ファイルのFHを、改めてファイルを管理する前記サーバー装置に送信する前記クライアント装置から、または、他の前記サーバー装置から、前記ファイルデータ記憶装置に記憶されているファイル、のFHを含むリクエストを受信すると、前記メタデータ記憶装置から、入力されたFHを識別子とするファイルのファイル詳細情報を取得して返信する、分散ファイルシステム制御方法。

#### 【請求項7】

前記サーバー装置が、さらに、

前記クライアント装置から、新たなディレクトリ、または、ファイルの第1の生成要求を受信すると、前記第1の生成要求に含まれる名前に依存しない方法で新たなディレクトリ、または、ファイルを管理する前記サーバー装置を選択し、

選択した前記サーバー装置に新たなディレクトリ、または、ファイルの第2の生成要求を送信して、生成された新たなディレクトリ、または、ファイルのFHを受信し、さらに、新たなファイルを生成した時はファイル詳細情報も受信し、

受信したFH、並びに、前記第1の生成要求に含まれる名前、および、ファイルタイプを前記メタデータ記憶装置に格納されているディレクトリに記憶し、さらに、受信した場合はファイル詳細情報を前記ファイル詳細情報表に記憶し、

また、他の前記サーバー装置から前記第2の生成要求を受信すると、a) ディレクトリを生成して前記メタデータ記憶装置に格納、または、b) ファイル詳細情報を生成して前記メタデータ記憶装置に格納すると共に前記第2の生成要求送信元に返信し、さらに、自装置のサーバーID及び自装置内一意の値から新たなディレクトリ、または、ファイルのFHを作成して、作成したFHを前記メタデータ記憶装置に格納すると共に前記第2の生成要求送信元に返信する、請求項6の分散ファイルシステム制御方法。

#### 【請求項8】

前記サーバー装置が、前記クライアント装置から、ディレクトリ階層における直下のディレクトリの名前を受信すると、前記メタデータ記憶装置に格納されているディレクトリ

10

20

30

40

50

から、入力された名前の F H とファイルタイプを取得して出力し、

前記クライアント装置が、まず、1) 上位のディレクトリを管理する前記サーバー装置に、直下のディレクトリの名前を入力して、直下のディレクトリの F H とファイルタイプを取得することを、前記パス名に名前が含まれるディレクトリについて、ルートディレクトリを管理する予め定められた前記サーバー装置を起点に、上位から順に繰り返すディレクトリ探索を実行して、パス名中の最下位ディレクトリの F H とファイルタイプを取得し、次に、2) 最下位ディレクトリの F H が包含するサーバーIDの前記サーバー装置に、2-1) ファイルの F H、および、ファイルタイプにたいする要求、並びに、2-2) ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストを送信して、a) 前記パス名のファイルの F H、ファイルタイプ、および、ファイル詳細情報を取得し、または、b) 前記パス名のファイルの F H、および、ファイルタイプを取得し、

最下位ディレクトリの F H が包含するサーバーIDの前記サーバー装置から、ファイル詳細情報を取得できなかったとき、当該サーバー装置から取得した前記パス名のファイルの F H を、改めてファイルを管理する前記サーバー装置に送信して、ファイル詳細情報を取得する、請求項 6 乃至請求項 7 の何れか 1 項の分散ファイルシステム制御方法。

【請求項 9】

前記クライアント装置が、

ディレクトリの、名前、F H、及び、ファイルタイプを関連付けてキャッシュするメタデータ記憶部を、さらに、備え、

前記ディレクトリ探索の過程で前記パス名から、ディレクトリの名前を取り出すと、a 1) 取り出した名前(以降、取得ディレクトリ名)が前記メタデータ記憶部にキャッシュされていないければ、前記取得ディレクトリ名のディレクトリの上位のディレクトリを管理する前記サーバー装置に前記取得ディレクトリ名を送信して、前記取得ディレクトリ名のディレクトリの F H とファイルタイプを得るとともに、前記メタデータ記憶部に、前記取得ディレクトリ名、並びに、前記サーバー装置から取得した F H とファイルタイプをキャッシュし、a 2) キャッシュされていれば前記メタデータ記憶部から前記取得ディレクトリ名のディレクトリの F H とファイルタイプを取得する、請求項 7 乃至請求項 8 の何れか 1 項の分散ファイルシステム制御方法。

【請求項 10】

ファイルを記憶するファイルデータ記憶装置と、

ディレクトリ階層における直下のディレクトリ、または、ファイルの、a) 識別子であって、当該ディレクトリ、または、当該ファイルを管理するコンピュータの識別情報であるサーバーIDを含むファイルハンドル(以降、F H と略記)、b) 名前、および、c) ファイルまたはディレクトリの区別を示すファイルタイプを関連付けて記憶するディレクトリ、および、前記ファイルデータ記憶装置に記憶されているファイルのファイル詳細情報を記憶するメタデータ記憶装置と、に接続され、ファイル詳細情報のキャッシュ用メモリ領域であるファイル詳細情報表を備える前記コンピュータに、

クライアント装置、若しくは、他の前記コンピュータから処理リクエストを受信、または、他の前記コンピュータにリクエストを送信するとともに、処理中または送信中のリクエスト数をカウントし、

クライアント装置から、ディレクトリ階層における直下のファイルの名前を含む、1) ファイルの F H、および、ファイルタイプにたいする要求、並びに、2) ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストを受信すると、a) 前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のファイルの F H とファイルタイプを取得し、b 1) カウントしたリクエスト数が所定閾値を超えておらず、c 1) 受信された名前のファイルのファイル詳細情報が前記ファイル詳細情報表にキャッシュされていれば、キャッシュされているファイル詳細情報を取得し、c 2) 前記ファイル詳細情報表にキャッシュされていないければ、取得したファイルの F H を含むリクエストを、当該 F H が包含するサーバーIDの前記コンピュータに送信してファイル詳細情報を取得して、前記ファイル詳細情報表にキャッシュし、取得したファイルの F H、ファイルタイプと

10

20

30

40

50

ファイル詳細情報とを前記クライアント装置に返信し、b2) カウントしたリクエスト数が所定閾値を超えている場合は、入力されたリクエストを変更してファイル詳細情報の出力を中止し、取得したファイルのFH、ファイルタイプを前記クライアント装置に返信し、

また、ファイルのFHを、改めてファイルを管理する前記コンピュータに送信する前記クライアント装置から、または、他の前記コンピュータから、前記ファイルデータ記憶装置に記憶されているファイル、のFHを含むリクエストを受信すると、前記メタデータ記憶装置から、入力されたFHを識別子とするファイルのファイル詳細情報を取得して返信する、処理を実行させるプログラム。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、サーバー装置、分散ファイルシステム、分散ファイルシステム制御方法、および、プログラム、特に、多数の計算ノードにより構成される並列計算機または分散処理環境において使用されるサーバー装置、分散ファイルシステム、分散ファイルシステム制御方法、および、プログラムに関する。

【背景技術】

【0002】

特許文献1には、ネットワークを介して接続された複数のストレージサーバー装置から構成される分散ファイル管理システムにおいて、ファイルやディレクトリを分散して管理する方法が記載されている。

20

【0003】

特許文献2には、階層ディレクトリ構造における各ディレクトリ、例えば、/ , usr, bin, tmp, xxx, yyy, zzzの情報を、ディレクトリ名称(識別名)によって決定される計算機に分散して保存する分散ファイルシステムのディレクトリ管理方法が記載されている。本文書の段落0019は、ディレクトリ情報を階層ディレクトリ構造と関係なくバラして各計算機に分散保存するから、アクセスが1つの計算機に集中せず負荷分散でき、システムの性能が低下しないとしている。

【0004】

特許文献3には、サーバーの負荷に応じてI/O (Input / Output) 処理プロセスを増減させる方法が記述されている。つまり、I/O要求が増加傾向か、減少傾向にあるかをその傾きから予測することによりプロセスの増減を実施する。

30

【0005】

非特許文献1には、メタデータサーバーを持たない分散ファイルシステムが記載されている。この分散ファイルシステムの各サーバーへの分散方法は、概ね以下のように動作する。まず、分散ハッシュテーブルのハッシュ値の範囲(例:  $0 \sim 2^{32} - 1$ )を各サーバーへ均等に割り当てる。次に、ファイルのオープンの際、与えられたファイルのパス名からハッシュ値を計算する。そして、そのハッシュ値を含む範囲に割り当てられたサーバーを目的のファイルが存在するサーバーであると決定する。

【先行技術文献】

40

【特許文献】

【0006】

【特許文献1】特許第5367470号公報

【特許文献2】特開平5-23341号公報

【特許文献3】特許第4089506号公報

【非特許文献】

【0007】

【非特許文献1】<http://gluster.readthedocs.org/en/latest/Quick-Start-Guide/Architecture/> [http://www.gluster.org/community/documentation/index.php/GlusterFS\\_Concepts](http://www.gluster.org/community/documentation/index.php/GlusterFS_Concepts)

50

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0008】

特許文献1が開示するファイルやディレクトリの分散管理方法には、以下の問題がある。

## 【0009】

第1に、キャッシュ蓄積部(メタデータキャッシュ)は、ストレージサーバ装置にしか存在しない(図22)。このため、指定されたパス名の中のすべてのディレクトリ、ファイルのメタデータがキャッシュされている場合でも、ユーザ端末装置(クライアント)とストレージサーバ装置の間で通信が最低1回は発生する。

10

## 【0010】

第2に、各ストレージサーバ装置は、各ユーザの利用できるディレクトリ配下についてそれぞれキャッシュを持つことになると考えられる。このため、ストレージサーバ装置のキャッシュ蓄積部には、多数のユーザが同時に利用するマルチユーザ環境において、ユーザ数を考慮した大容量の記憶域が必要になる。したがって、システム価格が増加する、あるいは、キャッシュの容量不足により期待した効果が得られない可能性が考えられる。

## 【0011】

第3に、パス名中にディレクトリが複数ある場合、ユーザ端末装置が応答を得るまでに複数回のストレージサーバ装置間の通信が必要になり、OneHop方式(段落0028)とは言えない。

20

## 【0012】

第4に、段落0021には、ユーザ利用ファイル名に分割されたデータファイルの番号が付加され、データファイルの番号に対応したデータファイルのファイル名を出力する旨、及び、ユーザ利用ファイルを複数のデータファイルに分割して複数のストレージサーバ装置において分散して記憶する旨が記載されている。この方法では、ユーザ利用ファイル名がユーザによって変更された場合、データファイルは各ストレージサーバ装置に分散されているため、新たな名前を反映させるために複数回のサーバ間通信が必要となると考えられる。

## 【0013】

30

第5に、段落0084以降の記載によると、read/writeのリクエストの処理の延長で、パス検索を行うと解釈できる。つまり、ファイルのオープンにおいて行われるパス検索の処理がread/write処理に含まれるため、その分I/O性能に影響する。また、read/writeシステムコールのインターフェースが規格と異なるため、POSIX(Portable Operating System Interface for UNIX)に準拠していないと考えられる。

## 【0014】

最後に、ユーザ端末装置がリクエストを送信する際、どのストレージサーバ装置へ送信するかをどのように決定するのかが不明瞭である。このため、以前に検索を実行しキャッシュが存在するストレージサーバ装置があったとしても、そのストレージサーバ装置以外のストレージサーバ装置へリクエストを送信し、最初からパス検索することが必要になる可能性があり、有効にキャッシュが機能するとは限らない。

40

## 【0015】

特許文献2が開示する分散ファイルシステムのディレクトリ管理方法は、ディレクトリ名称(識別名)によって保存先の計算機を決める。この方法では、段落0020のようなハッシュや名前の代わりにID(Identification)を利用する方法などを加味しても、分散のされ方に偏りが出る可能性を否定できない。特に、段落0015、0021のような名称の最初の一字で保存先の計算機を決める分散方法では、分散のされ方に偏りが出る可能性が高い。例えば、分散方法が、アルファベットa~mで始まる名称を有するディレクトリは計算機11aに配置するという場合、abc1, abc2, ..., abc1000という名称に同時にアクセスされると、すべて同一の計算機11aに要求が集中するという不都合が起きる。

50

## 【 0 0 1 6 】

また、ディレクトリ名称の変更において、不都合が生じたり、煩雑な処理が必要となったりすることが考えられる。

## 【 0 0 1 7 】

さらに、特許文献 2 の方法は、計算機が要求されたディレクトリ情報を記憶していない場合には、該ディレクトリ情報をサーバーに問い合わせる（段落 0 0 1 3、図 4）。この方法は、同一サーバーに処理が集中する可能性が有るだけでなく、ファイルシステムを構成するサーバーと計算機の合計台数が多くなりコスト高となる。

## 【 0 0 1 8 】

非特許文献 1 が開示するハッシュによる分散方法は、以下のような問題点がある。

10

## 【 0 0 1 9 】

まず、非特許文献 1 には、ファイル名が変更された場合、新たな名前に基づくハッシュ値から決まるサーバーにポインターファイルを置き、実際にファイルのデータが存在するサーバーを指し示すとある。つまり、ファイル名変更後は、該ファイルへのアクセス時に通常とは異なる処理が必要となり、その分余計なオーバーヘッドが生じることになる。ファイル名の変更は、エンドユーザが通常行う操作であり、処理効率低下、処理遅延が懸念される。

## 【 0 0 2 0 】

さらに、ハードリンク（異なる名前でも同一ファイルにアクセスするためのリンク機構）への対応が困難であり、また多数のファイルの生成において、ある範囲のハッシュ値が多数出現する可能性もある。

20

## 【 0 0 2 1 】

特許文献 3 が開示する I/O 処理プロセスの増減は、プロセスの生成、終了という比較的重いとされる処理を伴うため、この処理により CPU (Central Processing Unit) 資源を使用することで、本来のクライアントから要求された処理を邪魔してしまう可能性が有る。

## 【 0 0 2 2 】

上記の問題に対処するために、本発明にかかるファイルシステムは、ファイルのデータだけではなく、ファイルの位置情報などファイルシステムを管理するメタデータも複数のサーバー配下に分散させて配置することでメタデータサーバーを用いない構成とする。そのうえで、計算ノード（ファイルシステムのクライアント）が、複数のサーバーの中からアクセス対象のファイルが存在するサーバーを特定する為の手段を備える。

30

## 【 0 0 2 3 】

ただし、このシステムにはメタデータサーバーのようなサーバーが存在しないため、クライアントは特定のサーバーに問い合わせることはできない。このため、ファイルシステムの最上位のディレクトリを管理するルートサーバーを設定する。そして、例えばファイルのオープン処理（open システムコール）に際して、指定されたパス名に基づき、クライアント主導で、パス名中の各ディレクトリを管理するサーバーを、ルートサーバーから順に辿って特定することとした。

## 【 0 0 2 4 】

従って、本発明にかかるシステムにおいては、クライアントが、如何に効率的にアクセス対象のファイルやディレクトリが格納されているサーバーを特定できる手段を用意するかが最初の課題となる。特に、ディレクトリ名や、ファイル名を元にした（ハッシュなどの）方法によりサーバーを特定するのではなく、より偏りが起き難い方法で各サーバーへ分散させる場合であっても、効率的にアクセス対象のファイルやディレクトリが格納されているサーバーを特定することが課題となる。

40

## 【 0 0 2 5 】

一方、多数の計算ノードがあると、ファイルオープンの処理などにおいて、ほぼ同時にディレクトリ配下のファイルにアクセスが集中することで、特定のサーバーにリクエストが集中し、クライアントへの応答が遅延して TAT (Turn Around Time) が増大する可能性がある。したがって、多数のリクエストが特定のサーバーに集中した際の効率化がさらな

50



る課題となる。

【課題を解決するための手段】

【0026】

本発明にかかる分散ファイルシステム等は、ファイルを多数のクライアント装置がほぼ同時にオープンする場合、1)対象ファイルが存在するディレクトリを管理するサーバー装置へのオープン対象ファイルのファイルハンドルなどのリクエストと、対象ファイルを管理するサーバーへのファイル詳細情報のリクエストを1つにまとめる。これにより、サーバー装置が、両リクエストを同時に処理できるようにする。この際、2)ファイル詳細情報を対象ファイルが存在するディレクトリを管理するサーバーのメモリ上にキャッシュすることで、サーバー装置間の通信を削減する。

10

【0027】

さらに、同一ディレクトリ配下の異なるファイルのオープンの場合で、3)ほぼ同時に所定閾値以上のリクエストを受信した場合は、処理待ちが発生してTATが悪化することを防ぐため、上記1)のようにリクエストを1つにまとめて処理することを止める。この場合は、クライアント装置が、個別にファイル詳細情報のリクエストを、ファイルを管理するサーバー装置に送信することで効率化する。

【0028】

以上、1)乃至3)を反映して、本発明の実施の形態のサーバー装置等は、以下のよう

20

【0029】

本発明の1実施の形態のサーバー装置は、

ファイルを記憶するファイルデータ記憶装置と、ディレクトリ階層における直下のディレクトリ、または、ファイルの、a)識別子であって、当該ディレクトリ、または、当該ファイルを管理するサーバー装置の識別情報であるサーバーIDを含むファイルハンドル(以降、FHと略記)、b)名前、および、c)ファイルまたはディレクトリの区別を示すファイルタイプを関連付けて記憶するディレクトリ、および、前記ファイルデータ記憶装置に記憶されているファイルのファイル詳細情報を記憶するメタデータ記憶装置と、に接続され、

ファイル詳細情報のキャッシュ用メモリ領域であるファイル詳細情報表と、

クライアント装置、若しくは、他の前記サーバー装置から処理リクエストを受信、または、他の前記サーバー装置にリクエストを送信するとともに、処理中または送信中のリクエスト数をカウントする通信部と、

30

ファイルのFHを入力されると、FHが包含するサーバーIDの前記サーバー装置にファイル詳細情報を要求するリクエストを送信して、ファイル詳細情報を受信するファイル詳細情報問い合わせ部と、

クライアント装置から、ディレクトリ階層における直下のファイルの名前を含む、1)ファイルのFH、および、ファイルタイプにたいする要求、並びに、2)ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストが受信されたとき、a)前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のファイルのFHとファイルタイプを取得し、b1)前記通信部がカウントしたリクエスト数が所定閾値を超えておらず、c1)受信された名前のファイルのファイル詳細情報が前記ファイル詳細情報表にキャッシュされていれば、キャッシュされているファイル詳細情報を取得し、c2)前記ファイル詳細情報表にキャッシュされていなければ、取得したファイルのFHを前記ファイル詳細情報問い合わせ部に入力してファイル詳細情報を取得して、前記ファイル詳細情報表にキャッシュし、取得したファイルのFH、ファイルタイプとファイル詳細情報とを前記クライアント装置に返信し、b2)前記通信部がカウントしたリクエスト数が所定閾値を超えている場合は、入力されたリクエストを変更してファイル詳細情報の出力を中止し、取得したファイルのFH、ファイルタイプを前記クライアント装置に返信し、また、前記サーバー装置から返信されたファイルのFHを、改めてファイルを管理する前記サーバー装置に送信する前記クライアント装置から、または、他の前記サーバー装置の

40

50

前記ファイル詳細情報問い合わせ部から、前記ファイルデータ記憶装置に記憶されているファイル、のFHを含むリクエストが受信されると、前記メタデータ記憶装置から、入力されたFHを識別子とするファイルのファイル詳細情報を取得して返信するディレクトリエントリ参照/更新部と、を備える。

【0030】

また、本発明の1実施の形態の分散ファイルシステムの制御方法は、  
ファイルを記憶するファイルデータ記憶装置と、

ディレクトリ階層における直下のディレクトリ、または、ファイルの、a) 識別子であって、当該ディレクトリ、または、当該ファイルを管理するサーバー装置の識別情報であるサーバーIDを含むファイルハンドル(以降、FHと略記)、b) 名前、および、c) ファイルまたはディレクトリの区別を示すファイルタイプを関連付けて記憶するディレクトリ、および、前記ファイルデータ記憶装置に記憶されているファイルのファイル詳細情報を記憶するメタデータ記憶装置と、に接続され、ファイル詳細情報のキャッシュ用メモリ領域であるファイル詳細情報表を備える前記サーバー装置が、

クライアント装置、若しくは、他の前記サーバー装置から処理リクエストを受信、または、他の前記サーバー装置にリクエストを送信するとともに、処理中または送信中のリクエスト数をカウントし、

クライアント装置から、ディレクトリ階層における直下のファイルの名前を含む、1) ファイルのFH、および、ファイルタイプにたいする要求、並びに、2) ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストを受信すると、a) 前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のファイルのFHとファイルタイプを取得し、b1) カウントしたリクエスト数が所定閾値を超えておらず、c1) 受信された名前のファイルのファイル詳細情報が前記ファイル詳細情報表にキャッシュされていれば、キャッシュされているファイル詳細情報を取得し、c2) 前記ファイル詳細情報表にキャッシュされていなければ、取得したファイルのFHを含むリクエストを、当該FHが包含するサーバーIDの前記サーバー装置に送信してファイル詳細情報を取得して、前記ファイル詳細情報表にキャッシュし、取得したファイルのFH、ファイルタイプとファイル詳細情報とを前記クライアント装置に返信し、b2) カウントしたリクエスト数が所定閾値を超えている場合は、入力されたリクエストを変更してファイル詳細情報の出力を中止し、取得したファイルのFH、ファイルタイプを前記クライアント装置に返信し

、  
また、ファイルのFHを、改めてファイルを管理する前記サーバー装置に送信する前記クライアント装置から、または、他の前記サーバー装置から、前記ファイルデータ記憶装置に記憶されているファイル、のFHを含むリクエストを受信すると、前記メタデータ記憶装置から、入力されたFHを識別子とするファイルのファイル詳細情報を取得して返信する。

【0031】

また、本発明の1実施の形態のプログラムは、  
ファイルを記憶するファイルデータ記憶装置と、

ディレクトリ階層における直下のディレクトリ、または、ファイルの、a) 識別子であって、当該ディレクトリ、または、当該ファイルを管理するコンピュータの識別情報であるサーバーIDを含むファイルハンドル(以降、FHと略記)、b) 名前、および、c) ファイルまたはディレクトリの区別を示すファイルタイプを関連付けて記憶するディレクトリ、および、前記ファイルデータ記憶装置に記憶されているファイルのファイル詳細情報を記憶するメタデータ記憶装置と、に接続され、ファイル詳細情報のキャッシュ用メモリ領域であるファイル詳細情報表を備える前記コンピュータに、

クライアント装置、若しくは、他の前記コンピュータから処理リクエストを受信、または、他の前記コンピュータにリクエストを送信するとともに、処理中または送信中のリクエスト数をカウントし、

クライアント装置から、ディレクトリ階層における直下のファイルの名前を含む、1)

ファイルのFH、および、ファイルタイプにたいする要求、並びに、2)ファイル詳細情報にたいする要求の2つの要求をまとめたリクエストを受信すると、a)前記メタデータ記憶装置に格納されているディレクトリから、入力された名前のファイルのFHとファイルタイプを取得し、b1)カウントしたリクエスト数が所定閾値を超えておらず、c1)受信された名前のファイルのファイル詳細情報が前記ファイル詳細情報表にキャッシュされていれば、キャッシュされているファイル詳細情報を取得し、c2)前記ファイル詳細情報表にキャッシュされていなければ、取得したファイルのFHを含むリクエストを、当該FHが包含するサーバIDの前記コンピュータに送信してファイル詳細情報を取得して、前記ファイル詳細情報表にキャッシュし、取得したファイルのFH、ファイルタイプとファイル詳細情報とを前記クライアント装置に返信し、b2)カウントしたリクエスト数が所定閾値を超えている場合は、入力されたリクエストを変更してファイル詳細情報の出力を中止し、取得したファイルのFH、ファイルタイプを前記クライアント装置に返信し、

10

また、ファイルのFHを、改めてファイルを管理する前記コンピュータに送信する前記クライアント装置から、または、他の前記コンピュータから、前記ファイルデータ記憶装置に記憶されているファイル、のFHを含むリクエストを受信すると、前記メタデータ記憶装置から、入力されたFHを識別子とするファイルのファイル詳細情報を取得して返信する、処理を実行させる。

【発明の効果】

【0032】

20

本発明にかかるサーバ装置は、メタデータも複数のサーバ装置に分散配置されている分散ファイルシステムにおいて、クライアント装置による効率的なアクセス対象サーバ装置の特定、および、多数のリクエストが特定のサーバ装置に集中した際の効率的なアクセスを可能とする。

【図面の簡単な説明】

【0033】

【図1】図1は、第1の実施の形態にかかる分散ファイルシステム5の概要を示す構成図である。

【図2】図2は、第1の実施の形態にかかるクライアント装置、サーバ装置、および、外部記憶装置の内部構成を示す図である。

30

【図3】図3は、メタデータとサーバ装置との階層関係を表す概念図である。

【図4】図4は、ファイルのオープン、リードにおけるディレクトリ探索、リード要求の送信/応答の様子を示す図である(その1)。

【図5】図5は、ファイルのオープン、リードにおけるディレクトリ探索、リード要求の送信/応答の様子を示す図である(その2)。

【図6】図6は、ファイルのオープン、リードにおけるディレクトリ探索、リード要求の送信/応答の様子を示す図である(その3)。

【図7】図7は、クライアント装置側のファイルオープン処理の動作例のフローチャートである(その1)。

【図8】図8は、クライアント装置側のファイルオープン処理の動作例のフローチャートである(その2)。

40

【図9】図9は、サーバ装置側のファイルオープン処理の動作例のフローチャートである(その1)。

【図10】図10は、サーバ装置側のファイルオープン処理の動作例のフローチャートである(その2)。

【図11】図11は、サーバ装置側のファイルオープン処理の動作例のフローチャートである(その3)。

【図12】図13は、上位ディレクトリ管理サーバ装置側のファイル、ディレクトリ生成処理の動作例のフローチャートである。

【図13】図13は、生成サーバ装置側のファイル、ディレクトリ生成処理の動作例の

50

フローチャートである。

【図 1 4】図 1 4 は、クライアント装置側のマウント処理の動作例のフローチャートである。

【図 1 5】図 1 5 は、サーバー装置側のマウント処理の動作例のフローチャートである。

【図 1 6】図 1 6 は、コンピュータ装置の構成図である。

【発明を実施するための形態】

【0034】

< 第 1 の実施の形態 >

< 概要 >

図 1 は、本実施の形態にかかる分散ファイルシステム 5 の概要を示す構成図である。分散ファイルシステム 5 においては、図 1 が示すように、例えば複数の計算ノードであるクライアント装置 1（以降、クライアント 1 と略記）が相互結合ネットワーク 3 に接続されており、ユーザのジョブを実行する。また、例えば複数のサーバー装置 2（以降、サーバー 2 と略記）も同じ相互結合ネットワーク 3 に接続されており、任意のクライアント 1 とサーバー 2 の間、及び、任意のサーバー 2 同士の通信が可能となっている。各サーバー 2 の配下には少なくとも 1 台の外部記憶装置 4 が接続される。

10

【0035】

分散ファイルシステム 5 は、システム内のファイルを管理する為に、階層的なファイルディレクトリを用いる。

【0036】

20

仮に、分散ファイルシステム 5 に、メタデータサーバーと呼ばれるような、システム全体のメタデータを管理するサーバー 2 があると、メタデータを更新する処理、例えばファイルやディレクトリの生成や削除、が同時に多数発生した場合、メタデータサーバーにアクセスが集中し、分散ファイルシステム 5 のボトルネックになる可能性がある。この問題を避けるため、分散ファイルシステム 5 は、ファイルのデータだけではなく、ファイルの位置情報などファイルを管理するメタデータも複数のサーバー 2 に分散させて配置する。すなわち、分散ファイルシステム 5 は、メタデータサーバーを用いない構成を採用している。

【0037】

分散ファイルシステム 5 は、メタデータサーバーは用いないが、ファイルシステムのルートディレクトリを管理するサーバー 2（ルートサーバーと呼称）を用いる。ルートディレクトリとは、マウント対象のファイルシステムの先頭のディレクトリを指す。分散ファイルシステム 5 のシステム構築時にシステム管理者が、ルートサーバーを決めておく。

30

【0038】

なお、以降の説明において、あるディレクトリを管理するサーバー 2 は、当該ディレクトリを記憶し、クライアント 1 や他のサーバー 2 の要求に応じて当該ディレクトリに格納されている下位のディレクトリやファイルに関するデータを出力するサーバー 2 を指す。また、あるファイルを管理するサーバー 2 とは、当該ファイルおよび当該ファイルに関する詳細情報を記憶し、クライアント 1 や他のサーバー 2 の要求に応じて当該ファイルのデータを出力するサーバー 2 を指す。一台のサーバー 2 が、あるディレクトリを管理すると同時に、あるファイルを管理することも有る。

40

【0039】

ファイル、ディレクトリの生成の際に、上位ディレクトリを管理するサーバー 2 が、新たに生成するファイル、ディレクトリを管理するサーバー 2 を、後述するサーバー情報表 2 1 1 に登録されているサーバー 2 から選定する。この際のサーバー 2 の選定は、ファイルやディレクトリの名前、またはそれらに付される一意の ID (IDentification) などに依存しない方法を用いて行われる。この方法は、例えば、ラウンドロビンや一様乱数に基づく選定方法である。

【0040】

選定されたサーバー 2 の識別子であるサーバー ID は、選定されたサーバー 2 から選定元

50

のサーバー 2 に返される FH (File handle) に含まれる。FH は、生成されたファイル、ディレクトリの名称、ファイルタイプと共にメタデータとして、上位ディレクトリ毎に選定元のサーバー 2 が保持する。

【 0 0 4 1 】

ここで、FH は、ファイル、ディレクトリを分散ファイルシステム 5 内で一意に識別するための識別子であり、サーバー ID の他にサーバー 2 内で一意の数値も含む。また、ファイルタイプは、ファイルかディレクトリかの区別情報である。

【 0 0 4 2 】

この結果、各サーバー 2 は、自装置が管理する各ディレクトリ配下にあるファイル、ディレクトリの名とそれらの位置情報であるサーバー ID の一覧をディレクトリエントリとして持つことができる。

10

【 0 0 4 3 】

分散ファイルシステム 5 においては、各サーバー 2 がクライアント 1 からのリクエスト毎にそれぞれ独立してこのファイル、ディレクトリの生成処理を行うので、システム全体としての排他制御などは不要である。さらに、各サーバー 2 が、サーバー 2 の選定に際してラウンドロビンや一様乱数に基づく方法で行うので、ファイル、ディレクトリの名前、及びディレクトリ構造とは関係なく、ディレクトリ、ファイルをシステム内でほぼ均等に配置できる。例えば、サーバー 2 が 3 台ある場合、1 番目のサーバー 2 は 2、3、1 番目、2 番目のサーバー 2 は 3、1、2 番目、3 番目のサーバー 2 は 1、2、3 番目のサーバー 2 の順で割り当てる方法により、ディレクトリ、ファイルをシステム内でほぼ均等に配置できる。

20

【 0 0 4 4 】

ファイル名を変更する場合は、その上位ディレクトリのディレクトリエントリ内の名前を書き換えるだけでよい。さらに、ハードリンクは、例えばディレクトリエントリ内に名前は異なるが同一の FH、ファイルタイプを持つエントリを作り、対象となるファイルを管理するサーバー 2 で該ファイルのファイル詳細情報のリンク数に 1 を加算することで実現可能となる。

【 0 0 4 5 】

運用中に新たに、分散ファイルシステム 5 にサーバー 2 を追加した場合でも、既存のファイル、ディレクトリの配置は変わらないので、特別な処理は不要である。

30

【 0 0 4 6 】

なお、ルートサーバーは、ファイルシステムの先頭のディレクトリのみを管理するのではなく、他のサーバー 2 同様にサーバー 2 選定の対象にもなるため、データ、メタデータの管理に関して他のサーバー 2 との違いはない。

【 0 0 4 7 】

ファイル、ディレクトリの生成に際して選定されたサーバー 2 は、生成したファイル、ディレクトリの FH 毎のディレクトリエントリを設け、更新 / 参照時刻、パーミッションなどの詳細情報を保持する。さらにファイルを生成した場合は、ファイルのサイズ、チャンクサイズ、並列数などをファイル詳細情報に含めて記憶する。

【 0 0 4 8 】

これにより、分散ファイルシステム 5 の各サーバー 2 は、ルートサーバーを頂点とした、例えば、図 3 に示すような階層的なツリー構造を作る。

40

a) ファイルシステムの最上位のディレクトリを管理するルートサーバーは、自装置が記憶するルートディレクトリの中に、そのディレクトリ配下にある「子ディレクトリ」を管理するサーバー 2 を特定できる情報を記憶している。

b) 子ディレクトリを管理するサーバー 2 は、自装置が記憶する子ディレクトリの中に、そのディレクトリ配下にある「孫ディレクトリ」を管理するサーバー 2 を特定できる情報を記憶している。

c) 孫ディレクトリを管理するサーバー 2 は、自装置が記憶する孫ディレクトリの中に、そのディレクトリ配下にあるアクセス対象のファイルを管理するサーバー 2 を特定できる

50

情報を記憶している。

【 0 0 4 9 】

上記構造に基づき、クライアント 1 は、ファイルシステムをマウントする際、そのマウント先としてルートサーバーを指定し、ルートサーバーからルートディレクトリの FH とファイルシステムを構成している全サーバー 2 のサーバー ID と IP アドレスの対応表を受け取る。

【 0 0 5 0 】

次に、クライアント 1 は、ファイルのオープン処理などにおいてディレクトリ探索を行う。すなわち、クライアント 1 は、指定されたパス名の中の最上位ディレクトリの直下にあるディレクトリを管理するサーバー 2 をルートサーバーに問い合わせ、さらにその下のディレクトリを管理するサーバー 2 を問い合わせることを繰り返す。クライアント 1 は、このディレクトリ探索により、指定されたパス名のファイルを管理するサーバー 2 のサーバー ID を取得する。

10

【 0 0 5 1 】

図 4 乃至図 6 は、ファイルのオープン、リードにおけるディレクトリ探索、リード要求の送信 / 応答の様子を示す。

【 0 0 5 2 】

図 4 において、ユーザがクライアント 1 で実行する AP (Application Program) が、マウントポイントが "/mnt" であって、パスが "/mnt/home/user1/file1" であるファイルをオープンするシステムコールを発行したとする。クライアント 1 は、ルートサーバーにマウント時に受け取ったルートディレクトリの FH と共に配下のディレクトリ "home" を「名前」として送信して問い合わせ、ディレクトリ "home" の FH を受け取る (図 4 中の a)。さらに、クライアント 1 は、受け取った FH に含まれるサーバー ID から "home" を管理するサーバー 2 を特定してそのサーバー 2 にディレクトリ "user1" を「名前」として送信して問い合わせ、ディレクトリ "user1" の FH を受け取る (図中の b)。最後にクライアント 1 は、同様に "file1" を問い合わせ、ファイル "file1" の FH を受け取る (図 4 中の c)。

20

【 0 0 5 3 】

この結果、クライアント 1 は、AP が発行したファイル "file1" に対するリード要求を、当該ファイルを管理するサーバー # 4 に送信することが出来る (図 5 中の d)。

30

【 0 0 5 4 】

なお、この操作は、漸化式  $a_{i+1} = f(g(a_i), a_i, \text{名前})$  で表すことができる。ここで、 $a_i$  はパス名中の  $i$  番目の名前に対応する FH、初期値  $a_0$  はマウント時にルートサーバーから返されたルートディレクトリの FH である。さらに、 $g$  は FH に対応するディレクトリを管理するサーバー 2 のサーバー ID を返す関数、 $f$  はサーバー ID、FH、名前から、その名前に対応する FH を返す関数である。また、 $i$  は  $i = 0, \dots, n-1$  ( $n$  は、パス名中の名前の数) の整数である。

【 0 0 5 5 】

また、この過程において、クライアント 1 は、途中で受信したディレクトリを管理するサーバー 2 の情報をメタデータキャッシュ (図 4 乃至図 6 のクライアント 1 を参照) に保持する。このメタデータキャッシュは、後述するクライアント 1 のメタデータ記憶部 110 に格納される。そして後刻、上記のディレクトリ "user1" 配下の他のファイル、例えば "file2"、のオープン処理において、既に問い合わせ済みのディレクトリ名がパス名中に含まれる場合は、このメタデータキャッシュを参照することでサーバー 2 との通信回数を削減する。すなわちクライアント 1 は、サーバー # 1、# 2 との通信を省略し、サーバー # 3 からファイル "file2" の FH を受け取ることが出来る (図 6 中の a)。

40

【 0 0 5 6 】

なお、特許文献 2 は、頻繁にディレクトリ内容が書き換えられる場合、上記メタデータキャッシュについての不都合を指摘している (段落 0007、0008)。しかし、ファイルシステムの上位のディレクトリは、システム管理者でなければ書き換えはできないこ

50

とが通常であり、またエンドユーザが書き換え可能なディレクトリであっても、上位のディレクトリほど書き換えの頻度は少ない。このため、少なくとも上位のディレクトリについてはメタデータキャッシュが有効に機能する。

【0057】

また、ディレクトリ探索を実施する場合でも、その探索は比較的下位のディレクトリのみを対象とするため、ルートサーバーなどへ探索要求が集中することはなく、各サーバー2へ適度に負荷分散できる。

【0058】

次に、多数のクライアント1が同一ディレクトリ配下のファイルをオープンする場合の2つのケースについて、分散ファイルシステム5におけるI/Oの効率化方法を示す。

10

【0059】

ケース1：

同一ファイルへの並列I/Oなどの目的で、多数のクライアント1が同一ファイルをオープンする場合

ケース2：

多数のクライアント1が同一ディレクトリ配下の異なるファイルをオープンする場合

ここで、効率化の前提となる事項について説明する。

【0060】

ファイルのオープン処理において、ディレクトリ探索の最後の処理が、オープン対象のファイルのメタデータを取得することである。この際、クライアント1は、該ファイルのファイル詳細情報も必要としている。このファイル詳細情報は、例えば、ファイルのサイズ、更新/参照時刻、ファイルのチャンクサイズ、使用するサーバー数である。クライアント1は、この問い合わせを、該ファイルが存在するディレクトリを管理するサーバー2に対して行うが、これらファイル詳細情報は、もともと該ファイルを管理するサーバー2が保持している。

20

【0061】

このため、ファイル詳細情報をクライアント1が得る手段として下記2つの方法がある。

【0062】

(ア)該ファイルが存在するディレクトリを管理するサーバー2は、FHとファイルタイプのみをクライアント1へ返却する。その後、クライアント1は、返却されたFHから該ファイルを管理するサーバー2を特定し、そのサーバー2へファイル詳細情報を要求する。

30

【0063】

(イ)該ファイルが存在するディレクトリを管理するサーバー2が、該ファイルを管理するサーバー2へファイル詳細情報を問い合わせるサーバー2間通信(例えば、図6のa')を実施し、その結果をファイルのFHなどと共にクライアント1へ返す。つまり、本来は(ア)のように2つの要求であるものを1つの要求にまとめる。

【0064】

ここで、上述した2つのケースについて、I/Oの効率化方法を説明する。

【0065】

ケース1：

多数のクライアント1がほぼ同時期に上記(ア)によりファイル詳細情報を得る場合、クライアント1がサーバー2と通信する回数が、クライアント1毎に1回増える。クライアント1が512台あれば、クライアント1がサーバー2と通信する回数が512回となる。さらに、該ファイルを管理するサーバー2にもファイル詳細情報の要求が集中する可能性がある。ファイル詳細情報はデータ量が多少多くなるので、通信回数だけではなくサーバー2のI/OやCPU等の負荷の面からも好ましくない。

40

【0066】

このため、サーバー2間通信を行う(イ)の方法を基にして、該ファイルを管理するサーバー2から受け取った(例えば、図6中のa')ファイル詳細情報を、該ファイルが存

50

在するディレクトリを管理するサーバー 2 にキャッシュする。つまり、当該サーバー 2 は、該ファイルのFHと紐づくファイル詳細情報表 2 1 2 をキャッシュとしてメモリ上に持ち、ディレクトリエントリから参照可能とする。そして、当該サーバー 2 は、ファイル詳細情報表 2 1 2 にFHに対応するエントリがあれば、それを参照することで、ファイル詳細情報をクライアント 1 に出力する（例えば、図 6 中の a）。これにより、分散ファイルシステム 5 は、サーバー 2 間の通信回数（例えば、図 6 中のサーバー # 3 と # 5 の間）を削減し、さらに該ファイルを管理するサーバー 2（例えば、図 6 中のサーバー # 5）へのI/OやCPU等の負荷を減らす。

#### 【 0 0 6 7 】

##### ケース 2

対象ファイルがクライアント 1 毎にそれぞれ異なるため、(イ)の場合、サーバー 2 間通信は対象ファイルごとに必要になる。ただ、サーバー 2 ではケース 1 かケース 2 かの判別はできない。すなわち、次に来るリクエストが何であるかはそれを受け取るまでわからない。このため、サーバー 2 は、次に述べるデーモン数の範囲内では、ケース 1 と同様に(イ)の方法を選択する。

#### 【 0 0 6 8 】

(イ)の方法を使った場合、サーバー 2 間通信の間、ファイルが存在するディレクトリを管理するサーバー 2 上（例えば、図 6 中のサーバー # 3）のデーモンが対象ファイルごとに 1 つブロックされる。ここで、デーモンとは、クライアント 1 からのリクエストを処理するプロセスである。

#### 【 0 0 6 9 】

多数のクライアント 1 からのアクセスが集中し、デーモン数よりもクライアント 1 からのリクエスト数が多い場合は、デーモンが空くまでキューイングされて待たされることになり、その分TAT (Turn Around Time) が悪化することになる。ある特定のクライアント 1 のTATの悪化は、そのクライアント 1 を含むマルチノードジョブ全体のTATの悪化につながる。

#### 【 0 0 7 0 】

このため、サーバー 2 がリクエストを受信した際、デーモンがすべて処理中である場合は、サーバー 2 は、クライアント 1 のファイル詳細情報取得方法を、(イ)から(ア)の方法に切り替える。つまり、ファイルが存在するディレクトリを管理するサーバー 2 は、自装置が保持しているディレクトリエントリ内の情報であるFH、ファイルタイプのみをクライアント 1 に返す。クライアント 1 は、返されたFHから、オープン対象のファイルを管理するサーバー 2 を特定し、特定したサーバー 2 へ直接ファイル詳細情報を要求する。オープン対象の各ファイルは、それぞれ各サーバー 2 に分散配置されているため、このリクエストの送信先もクライアント 1 毎に分散されることになる。これにより、TATの悪化を軽減することが可能となる。

#### 【 0 0 7 1 】

##### < 構成 >

図 2 は、本実施の形態にかかるクライアント 1、サーバー 2、および、外部記憶装置 4 の内部構成を示す図である。

#### 【 0 0 7 2 】

クライアント 1 は、通信部 1 0 1、マウント要求部 1 0 2、メタデータ参照 / 更新部 1 0 3、サーバー情報表更新部 1 0 4、メタデータ問い合わせ部 1 0 5、パス名解析部 1 0 6、ファイル / ディレクトリ生成要求部 1 0 7、ファイル詳細情報要求部 1 0 8、サーバー情報表 1 0 9、および、メタデータ記憶部 1 1 0 を包含する。通信部 1 0 1 は、リクエスト作成部 1 0 1 - 1、および、サーバー特定部 1 0 1 - 2 を包含する。

#### 【 0 0 7 3 】

サーバー 2 は、通信部 2 0 1、マウント応答部 2 0 2、サーバー選定部 2 0 3、ファイル / ディレクトリ生成要求部 2 0 4、ディレクトリエントリ参照 / 更新部 2 0 5、ファイル / ディレクトリ生成部 2 0 6、FH生成部 2 0 7、サーバー内FH検索部 2 0 8、ファ

10

20

30

40

50



イル詳細情報問い合わせ部 209、I/O発行部 210、サーバー情報表 211、ファイル詳細情報表 212、および、メタデータ時記憶表 213を包含する。通信部 201は、リクエスト作成部 201-1、リクエスト応答部 201-2、リクエスト受信部 201-3、リクエスト内容変更部 201-4、リクエスト処理デーモン数判別/更新部 201-5、リクエスト処理デーモン数カウンタ 201-6、リクエスト内容判別部 201-7、および、サーバー特定部 201-8を包含する。

【0074】

外部記憶装置 4は、メタデータ記憶装置 401、および、ファイルデータ記憶装置 402を包含する。

【0075】

なお、図 2 中で各部を結ぶ矢印は、結ばれる両者間の主要な指示/情報の流れを示すものであるが、各部間の指示/情報の流れは、これらに限られるものではない。

【0076】

図 2 中の各部は、それぞれ概略次のように動作する。最初に、クライアント 1 に含まれる各部の概略動作について説明する。

【0077】

通信部 101は、送信デーモン、受信デーモンを包含する。そして、通信部 101は、サーバー 2 に対してリクエストを送信する際は、例えば、メタデータ問い合わせ部 105 から渡された要求を通信プロトコルにあった形式に変換して、リクエスト作成部 101-1 を用いてリクエストを作成する。また、この際に、通信部 101は、FHを基に、サーバー特定部 101-2 を用いて宛先のサーバー 2 の IP アドレスを得る。

【0078】

リクエスト作成部 101-1は、TCP/IPなど使用する通信プロトコルに沿った形式の packets を作成する。

【0079】

FHにはこれに対応するファイル、ディレクトリを管理するサーバー 2 のサーバー ID が含まれている。サーバー特定部 101-2は、FHを基にサーバー ID を抽出し、サーバー情報表 109を参照して、宛先サーバー 2 の IP アドレスを得る。

【0080】

マウント要求部 102は、クライアント 1 がファイルシステムを利用可能にするための処理を行う。マウントに際して、あらかじめシステム管理者が、ルートサーバーの IP アドレス若しくはマシン名、及び、マウントポイントを指定しておく必要がある。

【0081】

mount コマンドなどによりマウント処理が開始されると、クライアント 1 からルートサーバーにマウント要求が通信部 101 を介して送信される。その応答として、サーバー情報、および、ファイルシステムの先頭のディレクトリの FH がルートサーバーからクライアント 1 に返される。ここで、サーバー情報は、ファイルシステムを構成している全サーバー 2 のサーバー ID とその IP アドレスの対応表、および、各サーバー 2 配下のデバイス情報を包含する。

【0082】

サーバー情報は、クライアント 1 において、サーバー情報更新部 104 によりサーバー情報表 109 に記録され、FH はルートディレクトリの FH として、メタデータ参照/更新部 103 によりメタデータ記憶部 110 に記録される。

【0083】

メタデータ参照/更新部 103は、メタデータ記憶部 110 の参照と更新を行う。

【0084】

サーバー情報表更新部 104は、サーバー情報表 109 更新を行う。

【0085】

メタデータ問い合わせ部 105は、最初に、パス名解析部 106 により、指定されたパス名中のより上位の名前を 1 つ抽出する。例えば、マウントポイントが “/mnt”、パス

10

20

30

40

50

名が “ /mnt/home/user1/file1 ” である場合、まず “ home ” が抽出される。メタデータ問い合わせ部 1 0 5 は、それをメタデータ参照 / 更新部 1 0 3 により、ルートディレクトリ配下にこの名前がメタデータ記憶部 1 1 0 に既に登録されているかどうかを調べる。

【 0 0 8 6 】

記憶されていない場合、メタデータ問い合わせ部 1 0 5 は、マウント時にルートサーバーから取得したルートディレクトリのFHとその配下の名前 (“ home ” ) の組み合わせをルートサーバーに送信して、送信した名前に対応するFHを要求する。応答として、指定した名前に対応するFHとファイルタイプが返される。メタデータ問い合わせ部 1 0 5 は、これをメタデータ参照 / 更新部 1 0 3 により、メタデータ記憶部 1 1 0 に登録する。

【 0 0 8 7 】

その後、メタデータ問い合わせ部 1 0 5 は、パス名解析部 1 0 6 により、パス名中からその 1 つ下の名前 (“ user1 ” ) を抽出し、その上位のディレクトリ (“ home ” ) のFHからそれを管理するサーバー 2 を特定し、上記と同様にそのFHと名前 (“ user1 ” ) に対応するFHを要求する。メタデータ問い合わせ部 1 0 5 は、応答として返された名前に対応するFHとファイルタイプをメタデータ参照 / 更新部 1 0 3 により、メタデータ記憶部 1 1 0 に登録する。メタデータ問い合わせ部 1 0 5 は、同様に処理を繰り返す。

【 0 0 8 8 】

なお、パス名の末端の名前 (“ file1 ” ) は、オープンシステムコールから呼ばれている場合、ファイル名であるので、サーバー 2 からは名前に対応するFHと共にファイル詳細情報が返される ( 図 4 中の c ) 。

【 0 0 8 9 】

前述の図 4 は、メタデータ問い合わせ部 1 0 5 が関連モジュールを用いながら、ディレクトリ探索を経て、パス名に対応するファイルのFHとファイル詳細情報を取得する流れを示す。

【 0 0 9 0 】

パス名解析部 1 0 6 は、与えられたパス名から文字列操作により、名前を 1 つ抽出する。

【 0 0 9 1 】

ファイル / ディレクトリ生成要求部 1 0 7 は、与えられたパス名の最後の名前で、ファイルまたはディレクトリの作成を該当するサーバー 2 に要求する。与えられたパス名の最後の名前は、例えば、パス名が “ /mnt/home/user1/file1 ” なら、名前 “ file1 ” である。ただし、作成するファイル、ディレクトリの上位ディレクトリのFHがメタデータ記憶部 1 1 0 に登録されていない場合は、先にメタデータ問い合わせ部 1 0 5 により上位ディレクトリのFHの問い合わせを行う。

【 0 0 9 2 】

メタデータ問い合わせ部 1 0 5 の処理において、オープンシステムコールによりファイルをオープンする際、与えられたパス名の末端の名前に対応したFHの問い合わせでファイル詳細情報が返されず、FHとファイルタイプのみが返される場合がある。すなわち、図 4 中の c において、ファイル詳細情報が返されない場合がある。

【 0 0 9 3 】

ファイル詳細情報要求部 1 0 8 は、この場合、このFHに基づいて該当するサーバー 2 、例えば、図 4 のサーバー # 4 、へ直接ファイル詳細情報を要求し、それをメタデータ参照 / 更新部 1 0 3 によりメタデータ記憶部 1 1 0 に登録する。

【 0 0 9 4 】

サーバー情報表 1 0 9 は、ファイルシステムを構成している全サーバー 2 のサーバー ID とその IP アドレスの対応表と各サーバー 2 配下のデバイス情報をクライアント 1 に記憶するメモリ領域である。

【 0 0 9 5 】

メタデータ記憶部 1 1 0 は、メタデータをクライアント 1 にキャッシュするためのメモリ領域である。メタデータ記憶部 1 1 0 は、オープンシステムコールなどで指定されるバ

10

20

30

40

50

ス名中の各ディレクトリ名、またはファイル名をファイルディレクトリの上位から順に記憶し、それに対応するディレクトリ、ファイルのFHとファイルタイプを記憶する。さらに、ファイルの場合は、メタデータ記憶部 110 はファイル詳細情報も記憶する。

【0096】

クライアント装置 1 内の各部の機能分担は、適宜変更して実装しても良い。例えば、メタデータ問い合わせ部 105 は、パス名解析部 106、および、メタデータ参照/更新部 103 の機能を包含しても良い。

【0097】

次に、サーバー 2 に含まれる各部の概略動作について説明する。

【0098】

通信部 201 は、受信デーモンと送信デーモンとして動作し、リクエストの受信とその応答及びサーバー 2 間通信の際のリクエストの作成を行う。

【0099】

さらに、通信部 201 は、動作中/待ち状態のリクエスト処理デーモンの数を管理している。クライアント 1 からファイル詳細情報を要求された際、待ち状態のリクエスト処理デーモン数が 1 個以下の場合、通信部 201 はファイル詳細情報を返却せずに、当該ファイルのFHとファイルタイプを返却するようにリクエストの内容を変更する。なお、通信部 201 は、動作中/待ち状態のリクエスト処理デーモンの数を管理する代わりに、ファイル詳細情報取得のため入出力について、現在進行中の入出力数と発行可能な最大多重度と現在進行中の入出力数との差分、を管理しても良い。

【0100】

リクエスト作成部 201 - 1 は、サーバー 2 間通信の際に宛先のサーバー 2 を指定し、必要なパラメータの設定などをしてリクエストを作成し、送信デーモンに渡す。

【0101】

リクエスト応答部 201 - 2 は、クライアント 1、またはサーバー 2 間通信でのリクエストに対する応答に必要なデータを設定し、送信デーモンに渡す。

【0102】

リクエスト受信部 201 - 3 は、受信デーモンが受信したリクエストをリクエスト処理デーモンに渡す。なお、この際、リクエストがFH、ファイルタイプ、ファイル詳細情報を要求し、かつ、待ち状態のリクエスト処理デーモン数が 1 個以下であれば、リクエスト内容変更部 201 - 4 によりリクエストをFHとファイルタイプのみのリクエストに置き換える。

【0103】

リクエスト内容変更部 201 - 4 は、リクエスト内容の書き換えを行う。

【0104】

リクエスト処理デーモン数判別/更新部 201 - 5 は、処理中のリクエスト処理デーモンの数をリクエスト処理デーモン数カウンタ 201 - 6 に保持する。すなわち、リクエスト処理デーモン数判別/更新部 201 - 5 は、リクエスト処理デーモンの処理の開始/終了時に、リクエスト処理デーモン数カウンタ 201 - 6 のカウント値に 1 を加算/減算する。また、リクエスト処理デーモン数判別/更新部 201 - 5 は、受信デーモンがファイル詳細情報のリクエストを受信した際にリクエスト処理デーモン数カウンタ 201 - 6 のカウント値を参照する。

【0105】

リクエスト処理デーモン数カウンタ 201 - 6 は、処理中のリクエスト処理デーモンの数を持つカウンタであり、メモリ上にその領域を確保される。

【0106】

リクエスト内容判別部 201 - 7 は、クライアント 1 または他のサーバー 2 から受信したリクエストの種別を調べる。

【0107】

FHには、これに対応するファイル、ディレクトリを管理するサーバー 2 のサーバーIDが

10

20

30

40

50

含まれている。サーバー特定部 201 - 8 は、FHからサーバーIDを抽出し、サーバー情報表 211 を参照して宛先サーバー 2 のIPアドレスを特定する。

【0108】

マウント応答部 202 は、クライアント 1 から、ファイルシステムのマウントを要求するリクエストを受信する。そして、マウント応答部 202 は、同クライアント 1 に対し、サーバー情報表 211 に記録されている該ファイルシステムのルートディレクトリのFHと、該ファイルシステムを構成するサーバー 2 のサーバーIDとIPアドレスを返却する。さらに、マウント応答部 202 は、各サーバー 2 配下のデバイスに関する情報を返す。なお、同リクエストを受信したサーバー 2 がルートサーバーではない場合は、マウント応答部 202 は、同クライアント 1 へエラーを返す。

10

【0109】

サーバー選定部 203 は、ファイル、ディレクトリの生成を要求された際、サーバー情報表 211 を参照して、ファイル、ディレクトリの名前に依存しない方法、例えば、ラウンドロビンや一様乱数に基づき、サーバー 2 を一つ選定する。

【0110】

ファイル/ディレクトリ生成要求部 204 は、サーバー選定部 203 が選定したサーバー 2 に対してファイル、ディレクトリの生成を要求する。

【0111】

ディレクトリエントリ参照/更新部 205 は、ファイル、ディレクトリを生成した際、その上位ディレクトリを管理するサーバー 2 において、メタデータ記憶装置 401 内の上位ディレクトリのディレクトリエントリに、生成されたファイル、ディレクトリのエントリを追加する。また、ディレクトリエントリ参照/更新部 205 は、参照要求を受けて、ディレクトリエントリ内の指定された名前に一致するエントリ、あるいはディレクトリエントリの各エントリのFHとファイルタイプを返す。

20

【0112】

ファイル/ディレクトリ生成部 206 は、ファイル、ディレクトリの生成要求をサーバー 2 間通信により受け取って、メタデータ記憶装置 401 内に該当するエントリが有るかどうかをチェックする。無ければ、ファイル/ディレクトリ生成部 206 は、FH生成部 207 によりFHを新たに生成し、ファイルの場合はファイルの構成情報や、ファイルタイプと共に新たなエントリを記録する。

30

【0113】

FH生成部 207 は、ファイル、ディレクトリを生成する際に起動されて、当該サーバー 2 のサーバーIDとサーバー 2 内でユニークな数値を組み合わせ、ファイルシステム全体でユニークなFHを生成する。

【0114】

サーバー内FH検索部 208 は、I/O発行部 210 を介して、メタデータ記憶装置 401 内を検索し、指定されたFHを検索する。

【0115】

ファイル詳細情報問い合わせ部 209 は、クライアント 1 がオープンしようとする、自装置が管理するディレクトリ配下のファイルのファイル詳細情報を、当該ファイルを管理するサーバー 2 に問い合わせる。

40

【0116】

I/O発行部 210 は、メタデータ記憶装置 401、及びファイルデータ記憶装置 402 に対して、指定されたメモリ上のデータを書き込む、あるいは、指定されたデータをメモリ上に読み込む。この際、I/O発行部 210 は、メタデータの場合は同時にメタデータ時記憶表 213 にも登録し、さらに読み込みの際にはメタデータ時記憶表 213 に当該メタデータがあればその値を返す。

【0117】

サーバー情報表 211 は、ファイルシステムを構成するサーバー 2 のサーバーIDとIPアドレスの対応、及びルートディレクトリのFHを格納するメモリ上に確保された領域である

50

。また、メタデータ記憶装置 4 0 1 には、サーバー情報表 2 1 1 と同じデータが格納されており、サーバー 2 を起動させた際にサーバー情報表 2 1 1 に読み込まれる。

【 0 1 1 8 】

ファイル詳細情報表 2 1 2 は、ファイル詳細情報を格納するためのメモリ上に確保された領域である。

【 0 1 1 9 】

メタデータ時記憶表 2 1 3 は、メタデータをキャッシュするためのメモリ上に確保された領域である。

【 0 1 2 0 】

サーバー装置 1 内の各部の機能分担は、適宜変更して実装しても良い。例えば、ディレクトリエントリ参照 / 更新部 2 0 5 は、サーバー内 F H 検索部 2 0 8 の機能を包含しても良いし、ファイル / ディレクトリ生成部 2 0 6 は、F H 生成部 2 0 7 の機能を包含しても良い。

10

【 0 1 2 1 】

最後に、外部記憶装置 4 に含まれる各部が記憶する情報について説明する。

【 0 1 2 2 】

メタデータ記憶装置 4 0 1 は、ディレクトリエントリ、ファイル、ディレクトリの詳細情報、及びサーバー情報表 2 1 1 と同じサーバー 2 の構成情報などのファイルシステムの管理情報、いわゆるメタデータを記憶する記憶媒体である。

【 0 1 2 3 】

ファイルデータ記憶装置 4 0 2 は、ファイルのデータを記憶する記憶媒体である。なお、ファイルデータ記憶装置 4 0 2 とメタデータ記憶装置 4 0 1 は、物理的に別の媒体でも同一の媒体でも良い。

20

【 0 1 2 4 】

なお、クライアント 1 やサーバー 2 は、ファイル、ディレクトリの削除、属性情報の取得などのファイルシステムが通常持っている他の機能も備えているが、他の処理から容易に想到し得るため、ここでは説明を割愛する。

【 0 1 2 5 】

ここで、クライアント 1 における、通信部 1 0 1、マウント要求部 1 0 2、メタデータ参照 / 更新部 1 0 3、サーバー情報表更新部 1 0 4、メタデータ問い合わせ部 1 0 5、パス名解析部 1 0 6、ファイル / ディレクトリ生成要求部 1 0 7、および、ファイル詳細情報要求部 1 0 8 は、論理回路で構成される。各部は、適宜、クライアント 1 が備える図示されない半導体メモリにアクセスする。サーバー情報表 1 0 9、およびメタデータ記憶部 1 1 0 は、クライアント 1 が備える図示されない半導体メモリ上に設けられる。

30

【 0 1 2 6 】

また、サーバー 2 における、通信部 2 0 1、マウント応答部 2 0 2、サーバー選定部 2 0 3、ファイル / ディレクトリ生成要求部 2 0 4、ディレクトリエントリ参照 / 更新部 2 0 5、ファイル / ディレクトリ生成部 2 0 6、F H 生成部 2 0 7、サーバー内 F H 検索部 2 0 8、ファイル詳細情報問い合わせ部 2 0 9、および、I / O 発行部 2 1 0 は、論理回路で構成される。各部は、適宜、サーバー 2 が備える図示されない半導体メモリにアクセスする。サーバー情報表 2 1 1、ファイル詳細情報表 2 1 2、およびメタデータ時記憶表 2 1 3 は、サーバー 2 が備える図示されない半導体メモリに記憶される。

40

【 0 1 2 7 】

外部記憶装置 4 は、例えば、H D D (Hard-Disk Drive) や S S D (Solid State Drive) である。

【 0 1 2 8 】

また、クライアント 1、およびサーバー 2 は、それぞれ、プログラム 4 3 を備えるコンピュータ装置 4 0 で実現することも出来る。

【 0 1 2 9 】

図 1 6 は、コンピュータ装置 4 0 の構成図である。コンピュータ装置 4 0 は、バス 4 5

50

で相互に接続されたプロセッサ 4 1、主記憶部 4 2、外部記憶装置 4 4を備える。ここで、例えば、主記憶部 4 2は半導体記憶装置、外部記憶装置 4 4はHDDやSSDである。主記憶部 4 2はプログラム 4 3を記憶している。

【0130】

クライアント 1が記憶しているプログラム 4 3は、クライアント 1として用いられるコンピュータ装置 4 0において、プロセッサ 4 1で実行されることにより、プロセッサ 4 1を通信部 1 0 1、マウント要求部 1 0 2、メタデータ参照/更新部 1 0 3、サーバー情報表更新部 1 0 4、メタデータ問い合わせ部 1 0 5、パス名解析部 1 0 6、ファイル/ディレクトリ生成要求部 1 0 7、および、ファイル詳細情報要求部 1 0 8として機能させる。主記憶部 4 2は、サーバー情報表 1 0 9、および、メタデータ記憶部 1 1 0を格納する。

10

【0131】

さらに、サーバー 2が記憶しているプログラム 4 3は、サーバー 2として用いられるコンピュータ装置 4 0において、プロセッサ 4 1で実行されることにより、プロセッサ 4 1を通信部 2 0 1、マウント応答部 2 0 2、サーバー選定部 2 0 3、ファイル/ディレクトリ生成要求部 2 0 4、ディレクトリエントリ参照/更新部 2 0 5、ファイル/ディレクトリ生成部 2 0 6、FH生成部 2 0 7、サーバー内FH検索部 2 0 8、ファイル詳細情報問い合わせ部 2 0 9、および、I/O発行部 2 1 0として機能させる。主記憶部 4 2は、サーバー情報表 2 1 1、ファイル詳細情報表 2 1 2、およびメタデータ一時記憶表 2 1 3を格納する。

20

【0132】

サーバー 2において、外部記憶装置 4 4または主記憶部 4 2は、メタデータ記憶装置 4 0 1、および、ファイルデータ記憶装置 4 0 2として機能する。

【0133】

<動作>

次に、本発明の実施例の動作について詳細に説明する。なお、説明にあたって、クライアント 1、サーバー 2は、コンピュータ装置 4 0を用いて実現されていると仮定する。クライアント 1、サーバー 2のオペレーティングシステムは、UNIX(登録商標)やLinux(登録商標)であると仮定する。また、クライアント 1、サーバー 2は、EtherNet(登録商標)、InfiniBand(登録商標)など一般に利用可能なネットワークインターフェースを持ち、相互結合ネットワーク 3を介して接続される。これにより、クライアント 1は任意のサーバー 2と通信可能であり、サーバー 2はそのサーバー 2以外の任意のサーバー 2と通信可能である。

30

【0134】

1. ファイルのオープン処理

図 7乃至図 1 1は、ファイルオープン処理の動作例のフローチャートである。図 7及び図 8はクライアント 1側、図 9乃至図 1 1はサーバー 2側の動作フローチャートである。なお、これらのフローチャートは、クライアント 1が、ファイルディレクトリを検索して、オープン対象のファイルのFHとファイル詳細情報を取得するまでの流れを示しており、その後のオープン処理については割愛されている。その後のオープン処理は、公知技術で実現できる。

40

【0135】

1. a. ファイルのオープン処理(クライアント 1側)

クライアント 1の動作は、メタデータ問い合わせ部 1 0 5が中心となる。クライアント 1上で動作するAPから、分散ファイルシステム 5内のファイルに対してオープンシステムコールが呼び出されると、オペレーティングシステムのカーネルからメタデータ問い合わせ部 1 0 5が起動されて、図 7の処理が開始される。

【0136】

最初に、メタデータ問い合わせ部 1 0 5は、パス名解析部 1 0 6により、指定されたパス名中の上位の名前を 1つ抽出する(ステップ 1)。パス名の末端、すなわち文字列としての終端、もしくは、作成対象の名前を検出した場合(ステップ 2でYES)、メタデー

50

タ問い合わせ部 105 は、図 8 の処理に進む。

【0137】

パス名の末端、もしくは作成対象の名前以外を検出した場合（ステップ 2 の NO）、メタデータ問い合わせ部 105 は、抽出した名前がメタデータ記憶部 110 に既に登録されているかどうかを、メタデータ参照/更新部 103 により調べる（ステップ 3）。登録されている場合は（ステップ 4 で YES）、メタデータ問い合わせ部 105 はステップ 1 に戻る。

【0138】

登録されていない場合は（ステップ 4 で NO）、メタデータ問い合わせ部 105 は、まず、上位ディレクトリの FH から問い合わせ先のサーバー 2 を特定し（ステップ 5）、該サーバー 2 へ上位ディレクトリの FH、問い合わせ対象の名前を送信する（ステップ 6）。メタデータ問い合わせ部 105 は、問い合わせ先のサーバー 2 から FH とファイルタイプを受信すると、それらを名前と共にメタデータ参照/更新部 103 によりメタデータ記憶部 110 に登録する（ステップ 7）。

10

【0139】

図 7 のステップ 2 で YES となり、図 8 の処理に進んだ場合、メタデータ問い合わせ部 105 は、ファイル詳細情報要求部 108 により、オープン対象のファイルの名前のファイル詳細情報、及び FH を、そのファイルの上位ディレクトリを管理するサーバー 2 へ問い合わせる。

【0140】

20

図 8 において、メタデータ問い合わせ部 105 は、メタデータ参照/更新部 103 により、オープン対象のファイルの名前が既にメタデータ記憶部 110 に登録されているかを調べ（ステップ 11）、既に登録されている場合（ステップ 12 で YES）、動作を終了する。登録されていない場合（ステップ 12 で NO）、メタデータ問い合わせ部 105 は、上位ディレクトリの FH より問い合わせ先のサーバー 2 を特定し（ステップ 13）、オープン対象ファイルの名前に相当する FH、及びファイル詳細情報を含む問い合わせを該サーバー 2 へ送信する（ステップ 14）。なお、この問い合わせは、1) 上位ディレクトリを管理するサーバー 2 が保持するオープン対象ファイルの名前に相当する FH、及び、ファイルタイプの問い合わせと、2) オープン対象ファイルを管理するサーバー 2 が保持するファイル詳細情報の問い合わせの 2 つの問い合わせをまとめた問い合わせである。

30

【0141】

FH、ファイルタイプと共にファイル詳細情報を受信した場合（ステップ 15 で YES）、メタデータ問い合わせ部 105 は、ステップ 9 の実施後動作を終了する。ファイル詳細情報は返されず、FH とファイルタイプが返された場合（ステップ 15 で NO）、メタデータ問い合わせ部 105 は、返された FH とファイルタイプをメタデータ参照/更新部 103 により一旦メタデータ記憶部 110 に登録する（ステップ 16）。さらに、メタデータ問い合わせ部 105 は、ファイル詳細情報要求部 108 により、返された FH を基にオープン対象ファイルを管理するサーバー 2 を特定し（ステップ 17）、該サーバー 2 へ FH を指定してオープン対象のファイルのファイル詳細情報を問い合わせる（ステップ 18）。その後、メタデータ問い合わせ部 105 は、ファイル詳細情報要求部 108 が受信した該ファイル詳細情報を、FH、名前と共にメタデータ参照/更新部 103 によりメタデータ記憶部 110 に登録する（ステップ 19）。

40

【0142】

メタデータ問い合わせ部 105 による図 8 の動作が終了すると、クライアント 1 上では、メタデータ記憶部 110 上のファイルの名前、FH、ファイル詳細情報を用いて、オープンシステムコールの処理が続行される。

【0143】

1. b. ファイルのオープン処理（サーバー 2 側）

図 9 乃至 11 は、図 7 のステップ 6、図 8 のステップ 14、およびステップ 18 においてサーバー 2 への問い合わせをクライアント 1 が行った際のサーバー 2 側の処理の例を示

50

す。クライアント 1 からの当合わせはサーバー 2 において、通信部 201 により受信される。

【0144】

図 9 において、通信部 201 は、受信されたクライアント 1 からのリクエスト（問い合わせ）を、リクエスト内容判別部 201 - 7 により何のリクエストかを調べる。そのリクエストが FH で示されたディレクトリ配下のファイルのファイル詳細情報と FH、ファイルタイプを要求するものであった場合（ステップ 21 で YES）、通信部 201 は、リクエスト処理デーモン数判別 / 更新部 201 - 5 によりリクエスト処理デーモン数カウンタ 201 - 6 を参照し、待ち状態の同デーモンの数を確認する（ステップ 22）。

【0145】

同デーモン数が 2 個以上である場合（ステップ 23 で NO）、通信部 201 は、図 11 のファイルの FH、ファイルタイプ、及び、ファイル詳細情報取得処理に進む。1 個以下である場合（ステップ 23 で YES）、通信部 201 は、リクエスト内容変更部 201 - 4 により、リクエストの内容を該ファイルの FH とファイルタイプの問い合わせのみに書き換える（ステップ 24）。次に、通信部 201 は、そのリクエストをリクエスト処理デーモンに渡し（ステップ 25）、図 10 のファイルの FH とファイルタイプ取得処理に進む。

【0146】

なお、受信したリクエストが FH で示されたディレクトリ配下のファイルの FH とファイルタイプを要求するものであった場合（ステップ 21 で NO）、通信部 201 は、そのリクエストをリクエスト処理デーモンに渡し（ステップ 25）、図 10 のファイルの FH とファイルタイプ取得処理に進む。

【0147】

図 9 でステップ 25 を実施した場合は、図 10 において、ディレクトリエントリ参照 / 更新部 205 は、サーバー内 FH 検索部 208 により、受信した FH をキーとして該当するディレクトリを特定する（ステップ 31）。ディレクトリエントリ参照 / 更新部 205 は、さらにそのディレクトリから、名前をキーとして該当するエントリを特定することで、名前に対応する FH とファイルタイプを抽出し、要求元のクライアント 1 へ送信し（ステップ 32）、処理を終了する。

【0148】

図 9 で、リクエスト処理デーモン数が 2 個以上である場合（ステップ 23 で NO）、図 11 において、ディレクトリエントリ参照 / 更新部 205 は、サーバー内 FH 検索部 208 により、受信した FH をキーとして該当するディレクトリを特定する（ステップ 41）。ディレクトリエントリ参照 / 更新部 205 は、さらにそのディレクトリから、名前をキーとして該当するエントリを特定することで、名前に対応する FH とファイルタイプを抽出する（ステップ 42）。

【0149】

続いてディレクトリエントリ参照 / 更新部 205 は、ファイル詳細情報表 212 を検索し、該ファイルのファイル詳細情報が登録されているかどうかを確認する（ステップ 43）。

【0150】

登録されていた場合（ステップ 43 で YES）、ディレクトリエントリ参照 / 更新部 205 は、FH とファイルタイプと共に該ファイル詳細情報を要求元のクライアント 1 へ返却し（ステップ 44）、処理を終了する。登録されていない場合（ステップ 43 で NO）、該ファイルのファイル詳細情報をファイル詳細情報問い合わせ部 209 が、該ファイルを管理するサーバー 2 へ問い合わせる（ステップ 45）。このとき、サーバー特定部 201 - 8 が、該ファイルの FH からサーバー 2 を特定する。

【0151】

問い合わせ先の該サーバー 2、すなわち、該ファイルを管理するサーバー 2 からファイル詳細情報を受信すると、ディレクトリエントリ参照 / 更新部 205 は、これをファイル詳細情報表 212 へ登録し（ステップ 46）、要求元のクライアント 1 へ FH、ファイルタ

10

20

30

40

50



イブと共に該ファイル詳細情報を送信し（ステップ47）、処理を終了する。

【0152】

なお、問い合わせ先のサーバー2では、ディレクトリエントリ参照/更新部205が問い合わせ元のサーバー2（該ファイルが存在するディレクトリを管理するサーバー2）が送信したFH受信する。そして、ディレクトリエントリ参照/更新部205が、サーバー内FH検索部208により、メタデータ時記憶表213またはメタデータ記憶装置401を検索し、受信したFHに対応するファイル詳細情報を読み出して、問い合わせ元のサーバー2へ返却する。

【0153】

2. ファイル、ディレクトリの生成処理

10

図12および図13は、サーバー2におけるファイル、ディレクトリ生成処理の動作例のフローチャートである。

【0154】

クライアント1上で動作するAPから、create/openシステムコール、または、mkdirシステムコールが呼び出されると、オペレーティングシステムのカーネルからメタデータ問い合わせ部105が起動されて、当該システムコールにおいて指定されたパス名を基に、生成するファイル、ディレクトリを作成するディレクトリのFHを取得する。その後、ファイル/ディレクトリ生成要求部107が起動され、該当するサーバー2へその生成を要求する。

【0155】

20

図12は、ファイル、ディレクトリを作成するディレクトリを管理するサーバー2側、図13は、ファイル、ディレクトリを生成するサーバー2側の動作フローチャートである。ここで、ファイル、ディレクトリを作成するディレクトリとは、新たに生成されたファイル、ディレクトリの上位ディレクトリとなるディレクトリである。

【0156】

2. a. ファイル、ディレクトリの生成処理（上位ディレクトリを管理するサーバー2側）

図12において、クライアント1から生成要求を受信すると、ファイル/ディレクトリ生成要求部204が起動される。ファイル/ディレクトリ生成要求部204は、サーバー選定部203により新たに生成するファイル、ディレクトリを管理するサーバー2を選定し（ステップ51）、選定したサーバー2へファイル、ディレクトリの作成要求としてファイルタイプと名前、ファイル、ディレクトリのパーミッションに関する情報を送信する（ステップ52）。

30

【0157】

ファイル/ディレクトリ生成要求部204は、作成要求先のサーバー2から、生成したファイル、ディレクトリのFHを、ファイルの場合はファイル詳細情報も、受信する。ディレクトリエントリ参照/更新部205は、受信したFHとファイルタイプ及び名前をI/O発行部210を経てメタデータ時記憶表213、及びメタデータ記憶装置401へ新たなディレクトリエントリとして登録する（ステップ53）。

【0158】

40

また、ファイルを生成した場合は（ステップ54でYES）、ディレクトリエントリ参照/更新部205は、さらに、受信したファイル詳細情報をファイル詳細情報表212に登録する（ステップ55）。この後、通信部201が、受信したファイル、ディレクトリのFHを、ファイルを生成した場合はそのファイルのファイル詳細情報も、要求元のクライアント1へ送信して（ステップ56）、処理を終了する。

【0159】

2. b. ファイル、ディレクトリの生成処理（作成するサーバー2側）

図12のステップ52でファイル、ディレクトリの生成要求を受信したサーバー2は、図13のフローチャートのように動作する。

【0160】

50

生成要求を受信したファイル/ディレクトリ生成部206は、FH生成部207により自サーバー2のサーバーIDと自サーバー2内で一意な番号を組み合わせることで、ファイルシステム内で一意なFHを生成する(ステップ61)。次に、ファイル/ディレクトリ生成部206は、このFHと共にファイル、ディレクトリの詳細情報としてファイルのパーミッション、生成時刻などをI/O発行部210を経て、メタデータ時記憶表213及びメタデータ記憶装置401に登録する(ステップ62)。

【0161】

この後、ファイル/ディレクトリ生成部206は、要求元のサーバー2へこれら情報を返して(ステップ63)、処理を終了する。

【0162】

3. マウント処理

図14および図15は、マウント処理の動作例のフローチャートである。図14はクライアント1側、図15はサーバー2側の動作フローチャートである。

【0163】

3. a. マウント処理(クライアント1側)

図14において、クライアント1のユーザがmountコマンドを入力、または、クライアント1上のAPがマウントシステムコールを呼び出すと、マウント要求部102がオペレーティングシステムから起動される。mountコマンド、または、mountシステムコールは、ルートサーバーであるサーバー2を指定しており、マウント要求部102は、指定されているサーバー2へマウント要求を送信する(ステップ71)。

【0164】

マウント要求部102は、該サーバー2から、ルートディレクトリのFH、ファイルシステムを構成する全サーバー2のサーバーIDとIPアドレスの対応表を含む応答を受信し(ステップ72)、同対応表を、サーバー情報更新部104によりサーバー情報表109に登録する(ステップ73)。マウント要求部102は、ルートディレクトリのFHをメタデータ参照/更新部103によりメタデータ記憶部110に登録して(ステップ74)、処理を終了する。

【0165】

3. b. マウント処理(サーバ2側)

図15において、マウント応答部202は、サーバー情報表211を参照し、ファイルシステムを構成する全サーバー2のサーバーIDとIPアドレスの対応表、及びルートディレクトリのFHの値を読み出し(ステップ81)、要求元のクライアント1へ送信して(ステップ82)、処理を終了する。

【0166】

なお、サーバー2での各リクエストの処理において、リクエスト処理デーモン数判別/更新部201-5は、リクエスト処理デーモン数カウンタ201-6の値をリクエストの処理開始時に1増加させ、リクエストの処理終了時に1減少させる。

【0167】

<効果>

本実施の形態にかかるファイル分散システム5の効果は以下の通りである。

【0168】

第1の効果は、多数のクライアント1での同一ファイルへのオープン処理により、そのリクエストが対象ファイルを管理するサーバー2とその上位ディレクトリを管理するサーバー2に集中した場合において、サーバー2間の通信回数を削減できることである。

【0169】

その理由は、上位ディレクトリを管理するサーバー2が、当該ファイルの詳細情報を、ファイルを管理するサーバー2から最初に取得した際に、ファイル詳細情報表212に保持するからである。このため、2度目以降のリクエストにおいては、上位ディレクトリを管理するサーバー2とファイルを管理するサーバー2との間の通信が不要となる。

【0170】

10

20

30

40

50

第2の効果は、多数のクライアント1での同一ディレクトリ配下の異なるファイルへのオープン処理により、そのリクエストが当該ディレクトリを管理するサーバー2に集中した場合において、当該サーバー2がボトルネックになるのを防止できることである。

【0171】

その理由は、当該ディレクトリを管理するサーバー2のI/O多重度が高くなり、デーモンがすべて使用中になった場合、当該サーバー2がディレクトリエントリ内に保持しているFH、及びファイルタイプのみをクライアント1へ返すからである。その後そのFHを基にクライアント1が、オープン対象のファイルを管理するサーバー2を特定し、そのサーバー2に問い合わせを行う。オープン対象の各ファイルは、それぞれ各サーバー2に分散配置されているため、このリクエストの送信先もクライアント1毎に分散されることになる。このため、アクセスが集中するディレクトリを管理するサーバー2のデーモンがすべて使用中であったとしても待ちが発生せず、クライアント1に対するTATが悪化しない。

10

【0172】

第3の効果は、ファイルへのread/write処理だけではなく、ファイルのオープン処理などのメタデータアクセス処理も複数のサーバー2で分散処理できることである。つまり、多数のクライアント1が同一ファイルシステムにアクセスする場合でも、read/write処理、メタデータ処理の両面において負荷分散できる。

【0173】

その理由は、ファイルのデータだけではなく、ファイル/ディレクトリの詳細情報やディレクトリエントリ等のメタデータもファイル、ディレクトリ毎に複数のサーバー2に分散して配置することが可能だからである。

20

【0174】

第4の効果は、ファイル、ディレクトリの生成の際に行われる、生成対象ファイル、ディレクトリを管理するサーバー2の選定において、システム全体として管理する情報の更新や排他制御により、システム全体のスループットを妨げないことである。

【0175】

その理由は、その上位ディレクトリを管理するサーバー2が、生成対象ファイル、ディレクトリを管理するサーバー2を、サーバー情報表211から選定するからである。このため、サーバー2の選定に際して、システム全体として管理する情報の更新や排他制御が不要となる。

30

【0176】

第5の効果は、同一ディレクトリ配下に多数のファイル、ディレクトリを生成しても、特定のサーバー2配下にこれらファイルデータ、メタデータが偏って配置されることはなく、容量、負荷の両面において適正に分散させることが可能なことである。

【0177】

その理由は、上位ディレクトリを管理するサーバー2が、新たに生成するファイル、ディレクトリを管理するサーバー2をサーバー情報表211から選定する際、名前に依存しない、例えば、ラウンドロビンや一様乱数などに基づく方法を取るからである。

【0178】

第6の効果は、クライアント1が、2度目以降の同一ディレクトリ配下へのファイルのオープン処理などメタデータアクセスを伴う処理において、サーバー2との通信回数を削減できることである。

40

【0179】

その理由は、ファイルを管理するサーバー2を特定する際、クライアント1上にメタデータをキャッシュするからである。例えば、オープン処理で同じディレクトリ配下の複数のファイルを同一クライアント1が続けてアクセスすることは、決して少なくない。メタデータのキャッシュにより、ファイルアクセスの都度、ルートサーバーからファイルを管理するサーバー2を辿る処理は不要になる。このため、処理効率が大幅に向上する。

【0180】

第7の効果は、名前に関する問題が軽減され、エンドユーザが通常行うファイル名変更

50

やハードリンク利用の操作の処理が煩雑にならないことである。

【0181】

その理由は、ファイル名を変更する場合は、その上位ディレクトリのディレクトリエンタリ内の名前を書き換えるだけでよいからである。また、ハードリンクは、例えばディレクトリエントリ内に名前は異なるが同一のFH、ファイルタイプを持つエンタリを作り、対象となるファイルを管理するサーバー2で該ファイルのファイル詳細情報のリンク数をカウントアップすることで実現可能となる。このため、本実施の形態にかかるファイル分散システム5は、特許文献1、2、または、非特許文献1のシステム等において発生するような問題を生じない。

【0182】

以上、実施形態を参照して本願発明を説明したが、本願発明は上記実施形態に限定されるものではない。本願発明の構成や詳細には、本願発明のスコープ内で当業者が理解し得る様々な変更をすることができる。

【産業上の利用可能性】

【0183】

本発明は、HPC (High Performance Computing) やビッグデータ解析のような分野において利用できる。特に、単一のノード、単一のCPUでは現実的な時間での計算や解析が不可能な多量のデータや計算量の処理を、多数のノードに分割して行うような並列計算機システムまたは分散処理基盤におけるデータストアに利用できる。

【符号の説明】

【0184】

- 1 クライアント
- 1 クライアント装置
- 2 サーバー
- 2 サーバー装置
- 3 相互結合ネットワーク
- 4 外部記憶装置
- 5 分散ファイルシステム
- 40 コンピュータ装置
- 41 プロセッサ
- 42 主記憶部
- 43 プログラム
- 44 外部記憶装置
- 45 バス
- 101 通信部
- 101 - 1 リクエスト作成部
- 101 - 2 サーバー特定部
- 102 マウント要求部
- 103 メタデータ参照/更新部
- 104 サーバー情報表更新部
- 105 メタデータ問い合わせ部
- 106 パス名解析部
- 107 ファイル/ディレクトリ生成要求部
- 108 ファイル詳細情報要求部
- 109 サーバー情報表
- 110 メタデータ記憶部
- 201 通信部
- 201 - 1 リクエスト作成部
- 201 - 2 リクエスト応答部
- 201 - 3 リクエスト受信部

10

20

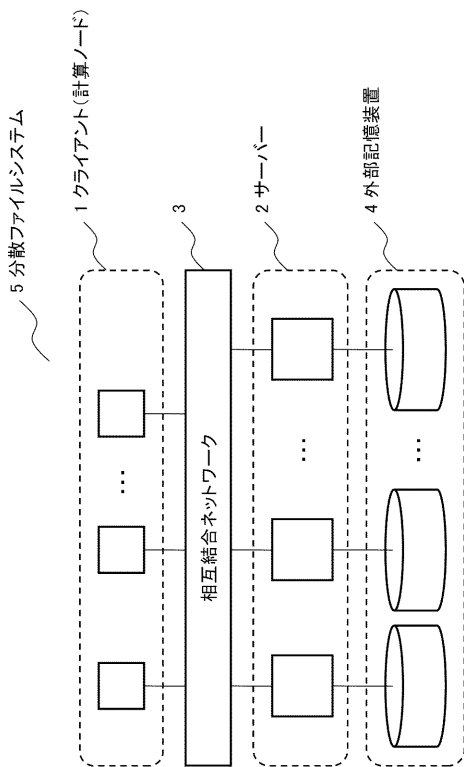
30

40

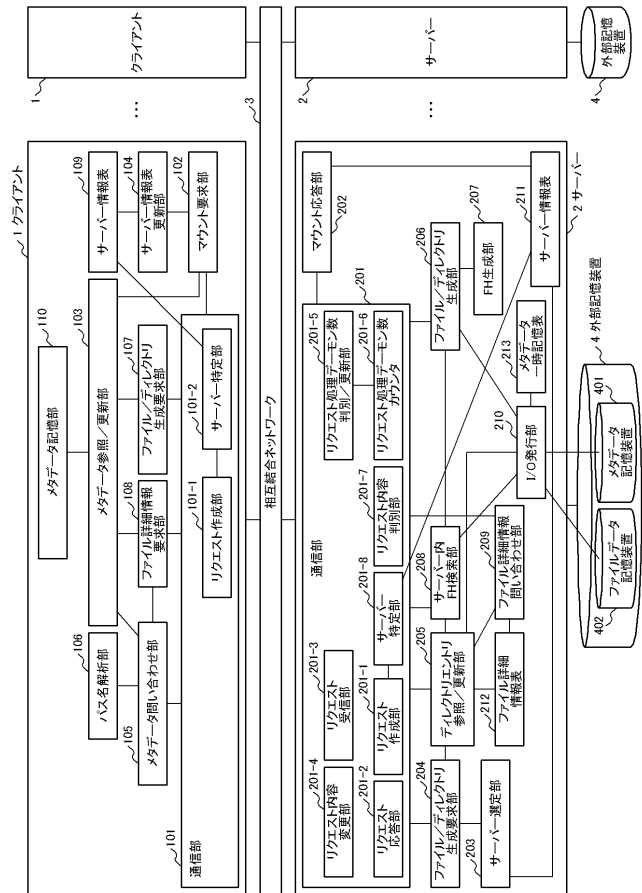
50

- 201 - 4 リクエスト内容変更部
- 201 - 5 リクエスト処理デーモン数判別/更新部
- 201 - 6 リクエスト処理デーモン数カウンタ
- 201 - 7 リクエスト内容判別部
- 201 - 8 サーバ特定部
- 202 マウント応答部
- 203 サーバ選定部
- 204 ファイル/ディレクトリ生成要求部
- 205 ディレクトリエントリ参照/更新部
- 206 ファイル/ディレクトリ生成部
- 207 FH生成部
- 208 サーバ内FH検索部
- 209 ファイル詳細情報問い合わせ部
- 210 I/O発行部
- 211 サーバ情報表
- 212 ファイル詳細情報表
- 213 メタデータ時記憶表
- 401 メタデータ記憶装置
- 402 ファイルデータ記憶装置

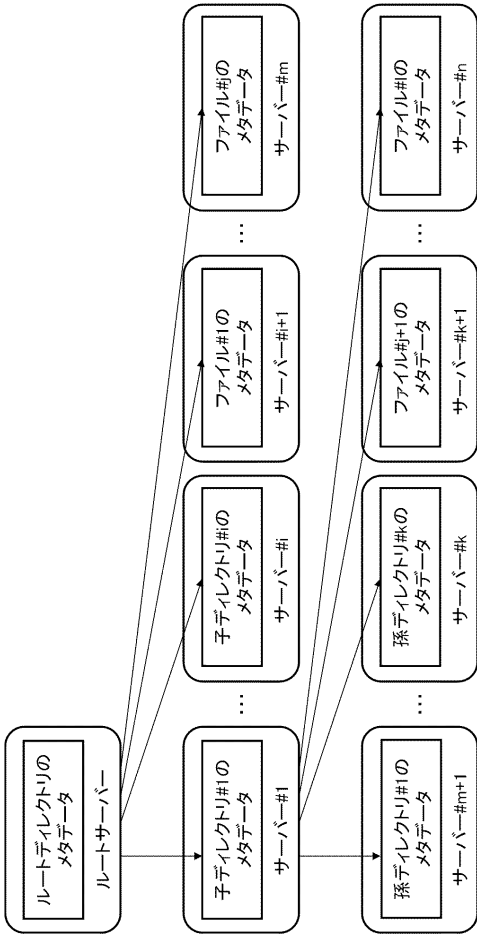
【 図 1 】



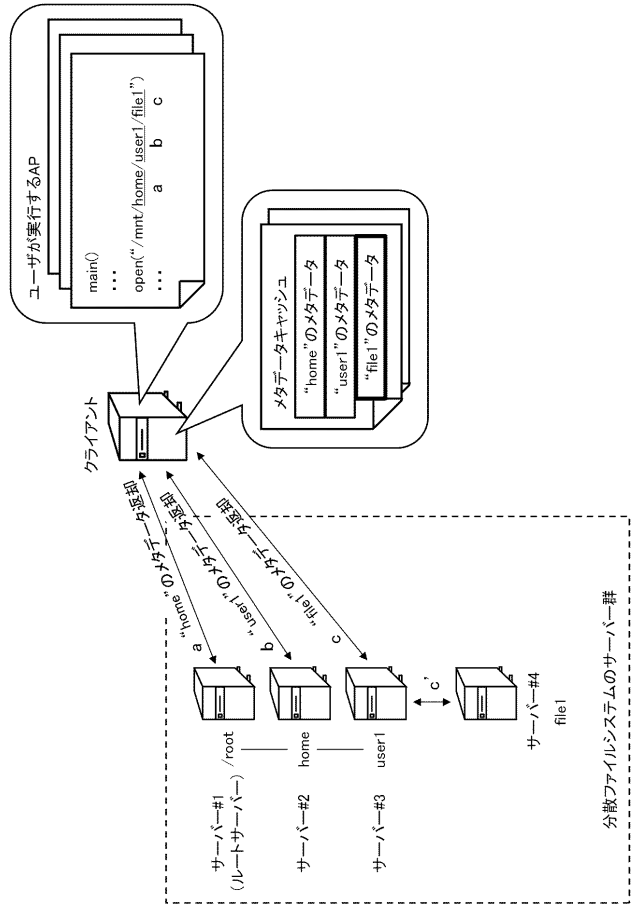
【 図 2 】



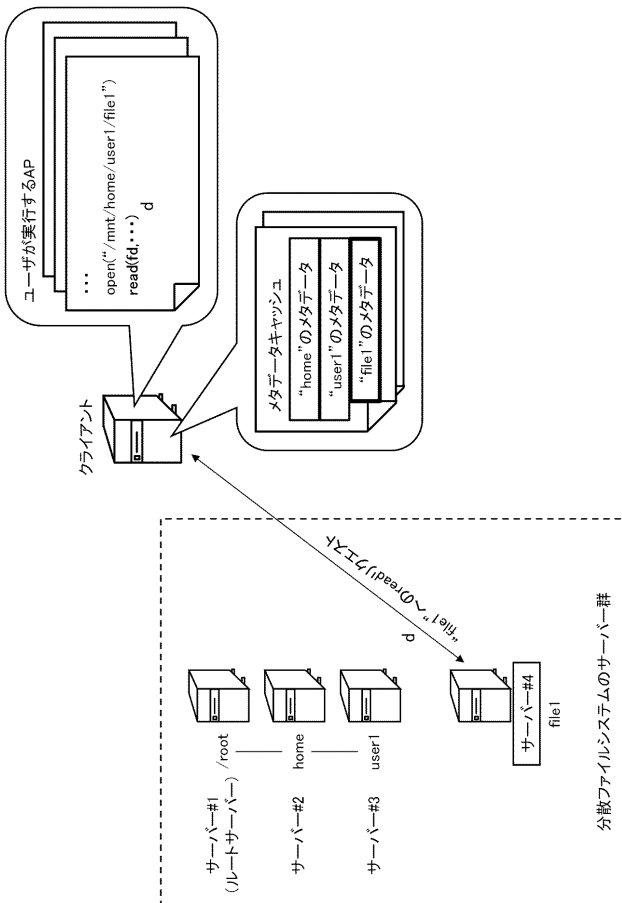
【 図 3 】



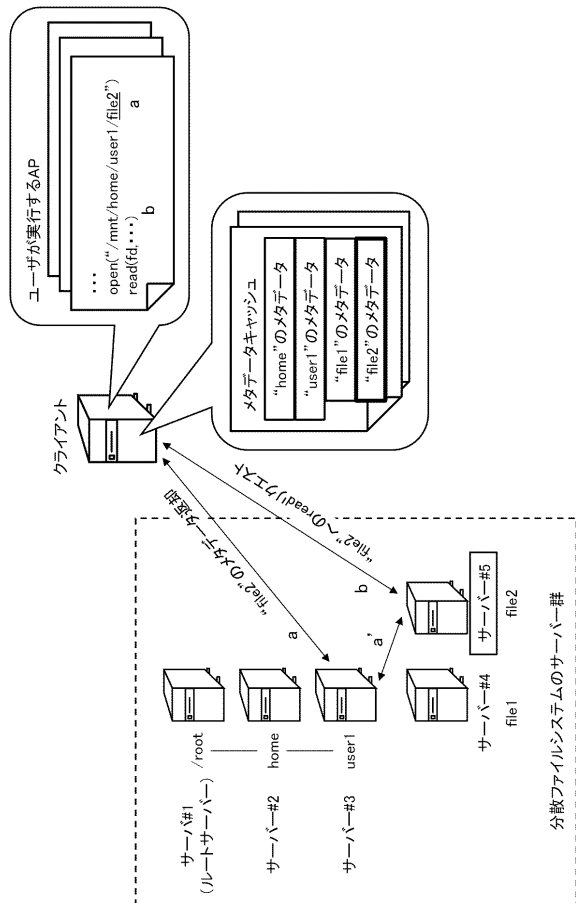
【 図 4 】



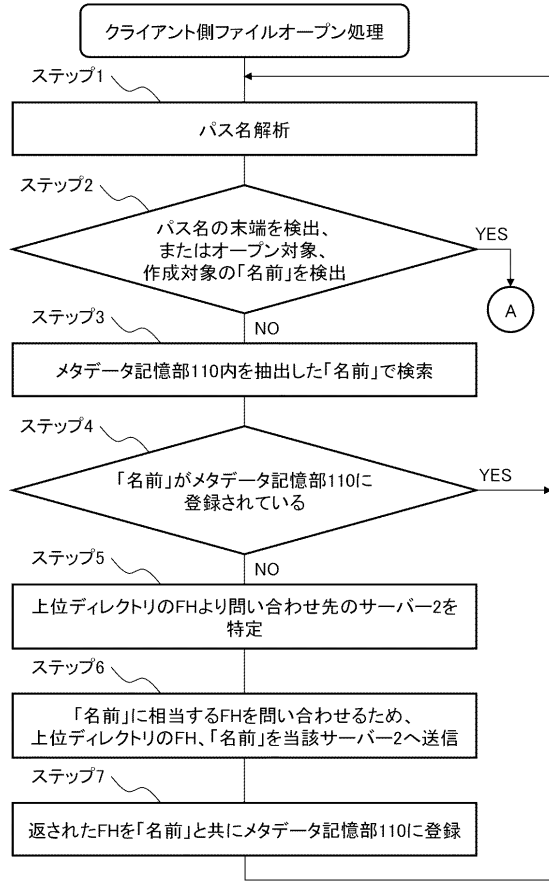
【 図 5 】



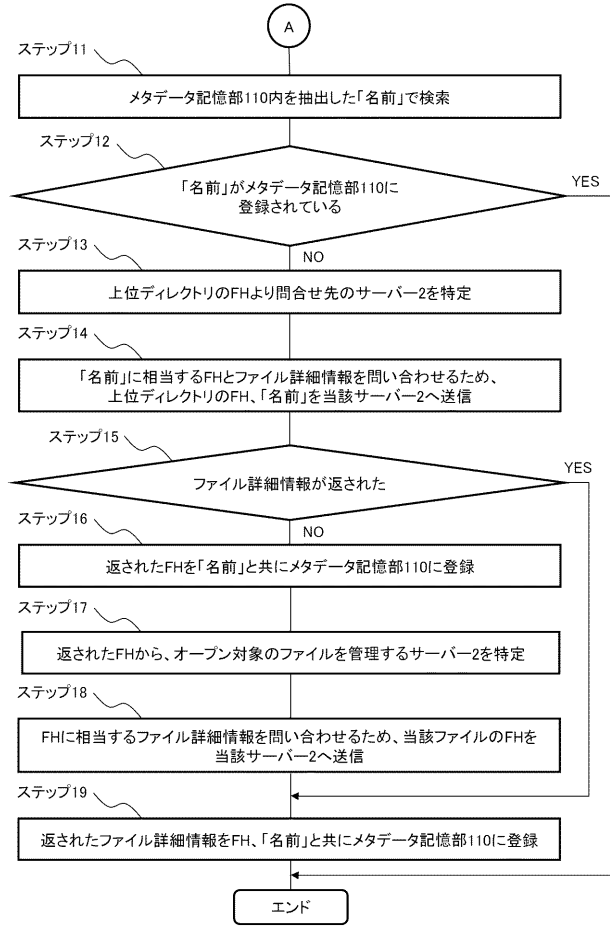
【 図 6 】



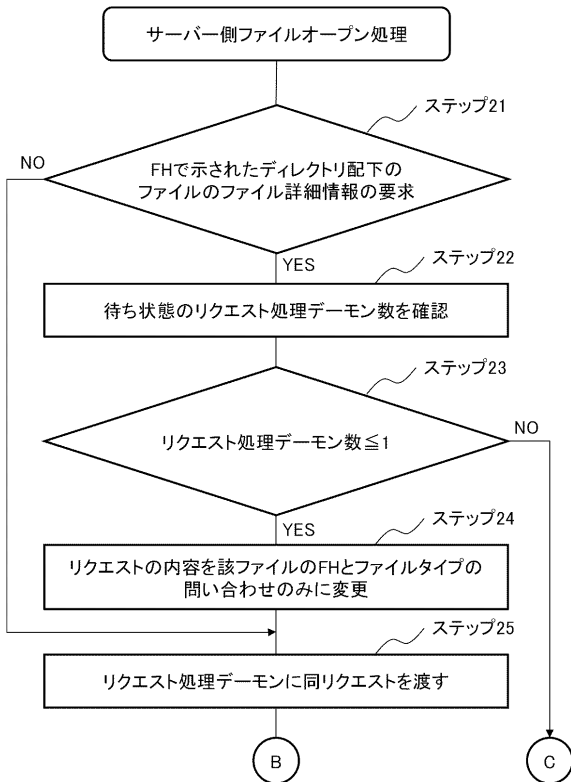
【 図 7 】



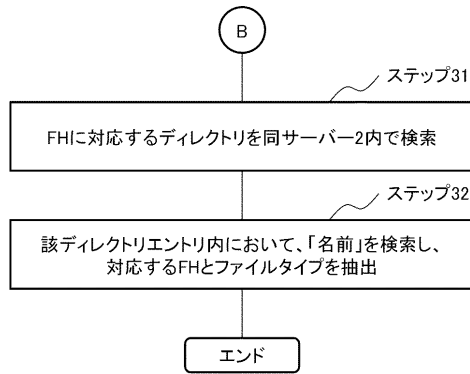
【 図 8 】



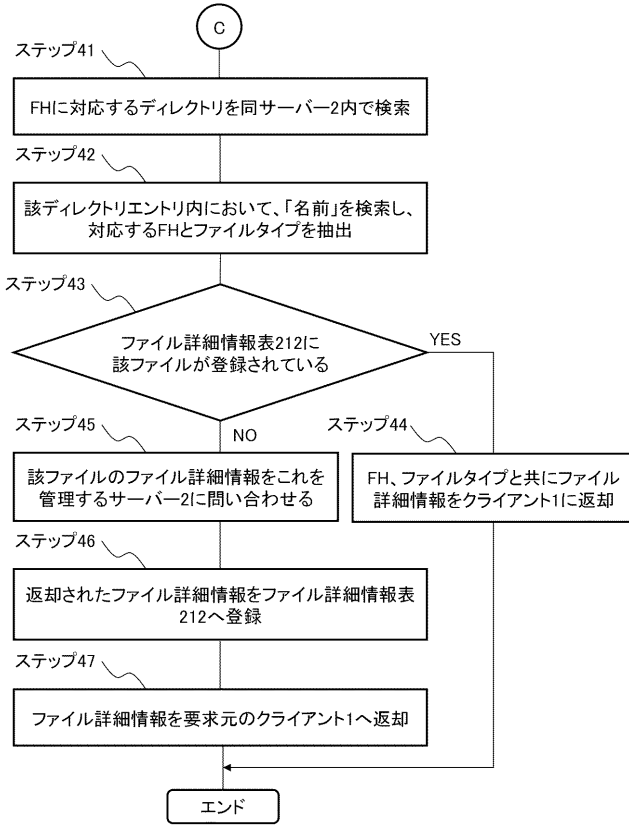
【 図 9 】



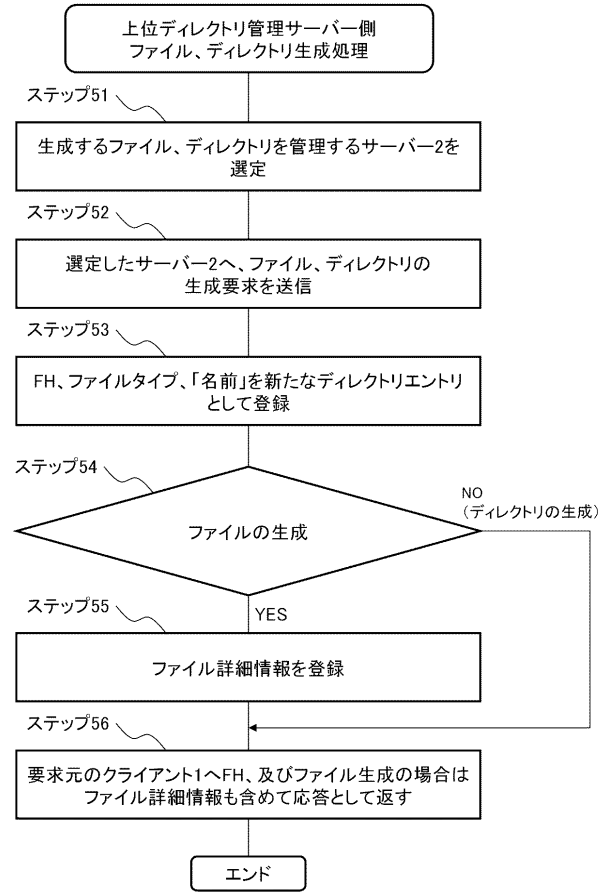
【 図 10 】



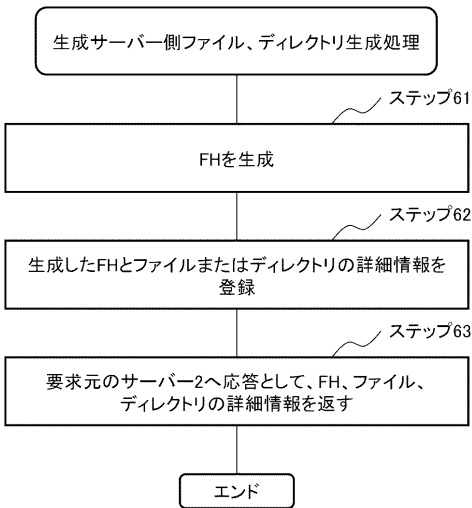
【 図 1 1 】



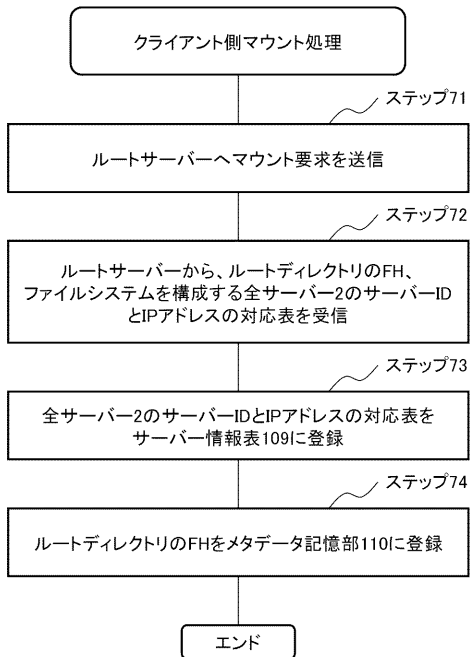
【 図 1 2 】



【 図 1 3 】

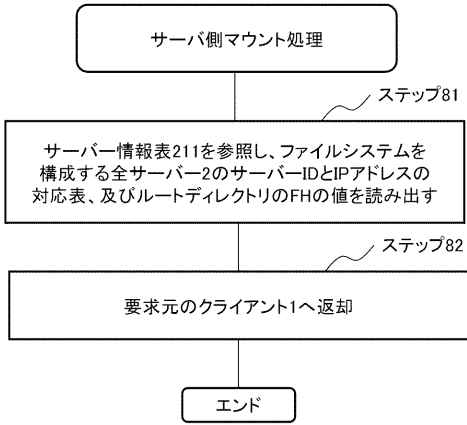


【 図 1 4 】





【図 15】



【図 16】

