

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6115396号
(P6115396)

(45) 発行日 平成29年4月19日(2017.4.19)

(24) 登録日 平成29年3月31日(2017.3.31)

(51) Int.Cl.

F I

G O 6 F 15/00 (2006.01)

G O 6 F 15/00 4 2 0 B

G O 6 F 12/00 (2006.01)

G O 6 F 12/00 5 3 5 B

G O 6 F 13/00 (2006.01)

G O 6 F 13/00 5 1 0 A

請求項の数 8 (全 34 頁)

(21) 出願番号 特願2013-168783 (P2013-168783)
 (22) 出願日 平成25年8月15日(2013.8.15)
 (65) 公開番号 特開2015-36928 (P2015-36928A)
 (43) 公開日 平成27年2月23日(2015.2.23)
 審査請求日 平成28年5月10日(2016.5.10)

(73) 特許権者 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番
 1号
 (74) 代理人 100092978
 弁理士 真田 有
 (74) 代理人 100112678
 弁理士 山本 雅久
 (72) 発明者 野崎 広弥
 神奈川県川崎市中原区上小田中4丁目1番
 1号 富士通株式会社内
 審査官 田中 幸雄

最終頁に続く

(54) 【発明の名称】 情報処理システム、情報処理装置、情報処理装置の制御プログラム、及び情報処理システムの制御方法

(57) 【特許請求の範囲】

【請求項 1】

情報処理装置と、前記情報処理装置との間で確立された接続を用いて前記情報処理装置と通信を行なう端末装置とを有する情報処理システムにおいて、

前記情報処理装置は、

前記接続を解除する予定時刻を前記端末装置に通知し、

前記端末装置は、

前記情報処理装置へ要求を送信する際に、現在時刻が前記情報処理装置から通知された予定時刻を経過しているか否かを判断する判断部と、

前記判断部により前記現在時刻が前記予定時刻を経過していると判断された場合、前記情報処理装置へ前記要求を送信する前に、前記情報処理装置との間で接続を確立するための接続要求を前記情報処理装置へ送信する送信部と、を有することを特徴とする、情報処理システム。

【請求項 2】

前記情報処理装置は、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断し、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除する接続解除部、を有する

ことを特徴とする、請求項 1 記載の情報処理システム。

【請求項 3】

10

20

前記情報処理装置は、
前記予定時刻を前記端末装置に関する接続情報に対応付けて管理する接続管理部、をさらに有し、
前記接続解除部は、
前記現在時刻が前記接続管理部の管理する前記予定時刻を経過していると判断した場合、前記接続情報を無効化する
ことを特徴とする、請求項 2 記載の情報処理システム。

【請求項 4】

前記情報処理装置は、
前記端末装置から接続要求又は前記接続要求以外の要求を受けると、前記接続要求又は前記要求への応答に前記予定時刻を含めて、前記端末装置へ通知する通知部、を有することを特徴とする、請求項 1 ~ 3 のいずれか 1 項記載の情報処理システム。

【請求項 5】

前記端末装置は、
前記情報処理装置から受信した予定時刻を前記情報処理装置に関する接続情報に対応付けて管理する時刻管理部と、
前記判断部により前記現在時刻が前記予定時刻を経過していると判断された場合、前記接続情報を無効化する接続情報管理部と、をさらに有する
ことを特徴とする、請求項 1 ~ 4 のいずれか 1 項記載の情報処理システム。

【請求項 6】

端末装置との間で確立した接続を用いて前記端末装置と通信を行なう情報処理装置において、
前記端末装置から接続要求を受けると、前記端末装置との間で接続を確立する接続処理を行なう接続管理部と、
前記接続処理により確立する接続を切り離す予定時刻を前記端末装置に通知する通知部と、
所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断し、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除する接続解除部、を有する
ことを特徴とする、情報処理装置。

【請求項 7】

端末装置との間で確立した接続を用いて前記端末装置と通信を行なう情報処理装置の制御プログラムにおいて、
前記情報処理装置に、
前記端末装置から接続要求を受けると、前記端末装置との間で接続を確立する接続処理を行なわせ、
前記接続処理により確立する接続を切り離す予定時刻を前記端末装置に通知させ、
所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断させ、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除させる
ことを特徴とする、情報処理装置の制御プログラム。

【請求項 8】

情報処理装置と、前記情報処理装置との間で確立された接続を用いて前記情報処理装置と通信を行なう端末装置とを有する情報処理システムの制御方法において、
前記情報処理装置が、
前記接続を解除する予定時刻を前記端末装置に通知し、
前記端末装置が、
前記情報処理装置へ要求を送信する際に、現在時刻が前記情報処理装置から通知された予定時刻を経過しているか否かを判断し、
前記現在時刻が前記予定時刻を経過していると判断した場合、前記情報処理装置へ前記

10

20

30

40

50

要求を送信する前に、前記情報処理装置との間で接続を確立するための接続要求を前記情報処理装置へ送信する

ことを特徴とする、情報処理システムの制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本件は、情報処理システム、情報処理装置、情報処理装置の制御プログラム、及び情報処理システムの制御方法に関する。

【背景技術】

【0002】

従来、分散ロック機能を有するクライアント・サーバモデルソフトウェア、例えばデータベースシステムや分散ファイルシステム（以下、これらをまとめて分散ファイルシステムという）が知られている。なお、分散ロック機能とは、クライアント・サービス（以下、単にクライアントという）からの共有リソースへのアクセスを制御する機能であり、例えば分散ロックマネージャ等の資源排他管理用サブシステムにより実現される。

【0003】

図15は、分散ロック動作の一例を示す図であり、図16は、クライアント追放時の分散ロック動作の一例を示す図である。なお、図16において図15と同一の符号が付された処理は、図15の処理と同様の処理であるため、重複した説明を省略する。

Lustre等の分散ファイルシステムでは、図15に示すように、サーバ・サービス（以下、単にサーバという）は、クライアントAからの分散ロック要求を受けると（処理T110）、クライアントAへ分散ロックを付与する（処理T120）。分散ロックが付与されたクライアントAは、ロック範囲へ例えば書込処理を行なう。

【0004】

また、サーバは、クライアントBから分散ロック要求を受け（処理T130）、クライアントAとクライアントBとの間で分散ロックの衝突を検出すると、クライアントAへ分散ロック（の返却）を要求する（処理T140）。しかし、クライアントAは、分散ロックを使用中であるため、分散ロック要求への返信では返却を拒否し（処理T150）、書込処理の完了後にサーバへ分散ロックを返却する（処理T160）。

【0005】

サーバは、返却された分散ロックをクライアントBへ付与する（処理T170）。処理T130から分散ロック付与を待つクライアントBは、付与された分散ロックのロック範囲へ例えば書込処理を行なう。なお、クライアントBは、サーバから分散ロック要求を受けていないため、書込処理完了後も分散ロックを返却しなくてよい。

以上のように、分散ファイルシステムの資源排他管理用サブシステムは、適宜適切な分散ロックを処理主体に付与し、分散ロックが付与されたクライアントまたはサーバのみが、資源を操作できるようにする（図15の「分散ロックの動き」参照）。これにより、複数のクライアントが同一ファイルに対して同時にwriteシステムコールを発行し、各クライアントが各々勝手にファイル書き込みを行なうことを防止できる。従って、同一データへの同時書き込みによるデータ破壊、データ喪失、管理情報の不整合によるファイルシステム破壊等の重大な障害の発生を防止できる。

【0006】

ここで、図16に示すように、サーバは、例えばクライアントからのリプライ（応答）がなかった場合（処理T250参照）等、クライアントの振る舞いが適切でなかった場合、当該クライアントの追放（追放処理）を行なう（処理T260参照）。追放処理では、サーバは、サーバから当該クライアントとの間の接続やロック範囲等に関する接続情報（サーバ側接続情報）を削除して、クライアントとの接続を切断する（解除する）。これにより、分散ファイルシステムでは、システム全体の整合性を維持することができる。

【0007】

ところで、クライアントの追放処理はサーバ主導の処理であり、クライアントに対して

10

20

30

40

50

非同期に実行される。また、仮にサーバがクライアントに対する通知・同期機能を有していたとしても、クライアントの追放処理は、クライアントがサーバに対して反応しない場合も想定した処理であるため、クライアントに対する通知が必ず成功することは保証されない。つまり、追放時の同期は保証されない。

【 0 0 0 8 】

図 1 7 は、クライアント追放処理及び追放復旧処理の一例を示す図である。

例えば、サーバ側で追放済（図 1 7 の処理 T 3 1 0 参照）のクライアントは、上述のように追放処理が非同期で行なわれるため、自身がサーバ側で追放されたことを認識できないことがある。従って、追放済のクライアントは、削除されたサーバ側接続情報とは不整合な状態である接続情報（クライアント側）に基づいて、サーバに対してリクエスト（要求）を送信する場合がある（処理 T 3 2 0）。当該リクエストを受信したサーバは、既に接続情報が存在しないことを確認すると、当該リクエストに対してエラーを返す（処理 T 3 3 0）。このとき、クライアントは自身がサーバ側で追放済であることを認識し、サーバ及びクライアント間でクライアントの追放処理に関する同期が完了する。

【 0 0 0 9 】

エラーを受けたクライアントは、不整合状態のクライアント側接続情報を全て破棄し、サーバ側に対して再接続要求を送信する（処理 T 3 4 0）。そして、クライアントは、再接続要求により確立された再接続によって接続情報を更新し、サーバと整合性のとれた接続情報を構築して、クライアントによる追放復旧処理が完了する（処理 T 3 5 0）。つまり、追放後、最初にクライアントからサーバへ送信したリクエストに対する、サーバからのエラーの通知が、クライアントによる追放復旧処理の契機となる。

【 0 0 1 0 】

ここで、上述したサーバからのエラーは、クライアント側接続情報とサーバ側接続情報との整合性がとれていないことに由来するエラーである。このエラーは、上述したように、データ破壊、データ喪失、管理情報の不整合によるファイルシステム破壊等の重大な障害をもたらす可能性を示唆するエラー（重大エラー）であるといえる。

図 1 8 は、クライアントにおけるエラー時のリトライ処理の一例を示す図であり、図 1 9 は、クライアント追放によるアプリケーションへの影響の一例を示す図である。図 1 8 に示すように、クライアントは、重大エラーを受けると、リトライ処理を行わずに、重大エラーを発生させたリクエストの発行契機となったおおもとの処理までエラーを返す。なお、図 1 8 に示すように、クライアントは、通常のエラーの場合、処理元までエラーを返さず、システム内部で折り返して、可能な限りリトライする。

【 0 0 1 1 】

例えば、L u s t r e 等において、重大エラーを発生させたリクエストの処理元がユーザ・アプリケーションの発行したシステムコールである場合（図 1 9 の処理 T 4 3 0 参照）、当該システムコールはエラー復帰することになる。つまり、ユーザ・アプリケーションに対して重大エラーが返されることになる（処理 T 4 4 0）。

また、重大なエラーを受けてクライアント側で追放復旧処理が開始されるが、直接的に追放復旧処理の契機を与えるリクエストを発行していない処理であっても、追放復旧処理前の古い接続情報を参照していた処理に対して、同じ重大エラーが返されることになる。当該処理がユーザ・アプリケーションに由来する処理であった場合には、やはりユーザ・アプリケーションに重大エラーが返されることになる。

【 0 0 1 2 】

つまり、クライアントの追放が実行されると、ユーザ・アプリケーションに対してエラーが返される可能性が高くなる。

ユーザ・アプリケーションは、ユーザが望む処理を行なうためのアプリケーションであり、多くの場合エラーを正しく処理することは考慮されていない。故に、ユーザ・アプリケーションは、重大エラーを受けた後もリトライするように作成されずに、エラー処理を正しく行なわない場合が多い。また、アプリケーションの実行が自動化されている場合等では、ユーザ・アプリケーションがエラー終了したことにユーザが暫く気付かないことも

10

20

30

40

50

多く、重大エラーの発生や取扱いがシステム運用上の障害となることも多い。

【0013】

そこで、L u s t r e 等の分散ファイルシステムでは、クライアントは、図20に示すように、サーバに対してpingリクエストを送信することがある。

図20は、pingによる追放検出手法の一例を示す図である。図20に示すように、クライアントは、サーバに対して定期的に（例えば25秒間隔で）pingリクエストを送信する（処理T510, T520）。これにより、クライアントは、pingリクエストがエラーになった契機（処理T520）で追放復旧処理を実行することができ（処理T530）、サーバ側で追放済のクライアントが長期間に亘って存在することを防止できる。また、ユーザ・アプリケーションが契機となって送信されたリクエスト（処理T540）が追放復旧処理の契機になる（重大エラーを発生させる）可能性を低減させることができる。

10

【0014】

この方法では、pingリクエストは、クライアントの状態及びサーバの状態と、完全非同期に、かつ定期的に実行される。従って、クライアント追放の観点では、pingリクエストの送信間隔ごとに、クライアント側接続情報及びサーバ側接続情報を同期させることが可能となる。

なお、関連する技術として、特定地域情報にアクセスを求める遠隔地からのユーザによって開始された要求に応答する技術が知られている（例えば、特許文献1参照）。この技術では、インターネットに接続し、動的に割り当てられるIP（Internet Protocol）アドレスを受け取り、当該IPアドレスを転送し、最大の使用されていない時間を超過したとき、接続解除を行なう。

20

【0015】

また、関連する他の技術として、クライアントから一定時間要求がなかった場合、ロックしている記憶装置内の領域の解放要求を記憶装置に対して発行するキャッシュストレージ装置が知られている（例えば、特許文献2参照）。

【先行技術文献】

【特許文献】

【0016】

【特許文献1】特表2003-521765号公報

30

【特許文献2】特開2004-342071号公報

【発明の概要】

【発明が解決しようとする課題】

【0017】

上述のように、クライアント追放処理は、サーバ主導でクライアントとは非同期的に実行される。つまり、クライアントは、実際にサーバに対してリクエストを送信し、リクエストに対する、クライアントの追放に起因するエラーをサーバから受信するまで、自身が追放されているか否かを認識できず、追放復旧処理を実行することが困難であった。従って、エラーを発生させたリクエストの送信元がユーザ・アプリケーションであった場合、ユーザ・アプリケーションにエラーが返ってしまうという課題がある。

40

【0018】

なお、図20に示す手法では、クライアントは、サーバに対して定期的にpingリクエストを送信し、pingリクエストがクライアントの追放に起因するエラーになったことを検知することで、追放復旧処理を行なう。しかし、システム規模が拡大し、サーバ数やクライアント数が増大するに従って、定期的にpingリクエストを送信することによるサーバ及びクライアントのCPU（Central Processing Unit）負荷やメモリ使用量、ネットワーク負荷等が膨大になる。

【0019】

また、負荷低減のために、pingリクエストを止めることも考えられる。しかし、追放処理後に最初に送信されるリクエストがエラーを発生させてしまい、当該リクエストの

50

送信元がユーザ・アプリケーションであった場合、ユーザ・アプリケーションにエラーが返ってしまうことになる。

なお、上述した関連する技術では、上述した課題については考慮されていない。

【0020】

ここまで、Lustre等の分散ファイルシステムを例に挙げて説明したが、上述した課題は、情報処理装置による端末装置の追放（接続の解除）処理が端末装置と非同期で行なわれる、種々の情報処理システムにおいても同様に生じ得る。

1つの側面では、本発明は、アプリケーションに対してエラーが返される頻度を低減させ、又は、情報処理システムの負荷の増大を抑制させて、情報処理装置と端末装置との間の接続を再確立することを目的とする。

【0021】

なお、前記目的に限らず、後述する発明を実施するための形態に示す各構成により導かれる作用効果であって、従来の技術によっては得られない作用効果を奏することも本発明の他の目的の1つとして位置付けることができる。

【課題を解決するための手段】

【0022】

本件の情報処理システムは、情報処理装置と、前記情報処理装置との間で確立された接続を用いて前記情報処理装置と通信を行なう端末装置とを有する情報処理システムにおいて、前記情報処理装置は、前記接続を解除する予定時刻を前記端末装置に通知し、前記端末装置は、前記情報処理装置へ要求を送信する際に、現在時刻が前記情報処理装置から通知された予定時刻を経過しているか否かを判断する判断部と、前記判断部により前記現在時刻が前記予定時刻を経過していると判断された場合、前記情報処理装置へ前記要求を送信する前に、前記情報処理装置との間で接続を確立するための接続要求を前記情報処理装置へ送信する送信部と、を有する。

【発明の効果】

【0023】

一実施形態によれば、アプリケーションに対してエラーが返される頻度を低減させ、又は、情報処理システムの負荷の増大を抑制させて、情報処理装置と端末装置との間の接続を再確立することができる。

【図面の簡単な説明】

【0024】

【図1】一実施形態の一例としての分散ファイルシステムの構成例を示す図である。

【図2】図1に示すFSクライアントに着目した分散ファイルシステムの構成例を示す図である。

【図3】図1に示すサーバ及びFSクライアントそれぞれのハードウェア構成例を示す図である。

【図4】図1に示すサーバの機能構成例を示す図である。

【図5】一実施形態の一例としてのサーバ及びクライアントが保持する追放予定時刻を説明する図である。

【図6】図4に示すリクエスト処理部による追放予定時刻の通知手法の一例を示す図である。

【図7】図2に示すクライアントの機能構成例を示す図である。

【図8】図2に示すサーバ及びクライアント間の通信の一例を説明するシーケンス図である。

【図9】図2に示すサーバによるリクエスト受信処理の一例を説明するフローチャートである。

【図10】図2に示すサーバによるクライアント追放処理の一例を説明するフローチャートである。

【図11】図2に示すクライアントによる接続リクエスト送信処理の一例を説明するフローチャートである。

10

20

30

40

50

【図 1 2】図 2 に示すクライアントによるリクエスト送信処理の一例を説明するフローチャートである。

【図 1 3】図 2 に示すクライアントによる返信受信待ち処理の一例を説明するフローチャートである。

【図 1 4】図 2 に示すクライアントによる返信受信処理の一例を説明するフローチャートである。

【図 1 5】分散ロック動作の一例を示す図である。

【図 1 6】クライアント追放時の分散ロック動作の一例を示す図である。

【図 1 7】クライアント追放処理及び追放復旧処理の一例を示す図である。

【図 1 8】クライアントにおけるエラー時のリトライ処理の一例を示す図である。

10

【図 1 9】クライアント追放によるアプリケーションへの影響の一例を示す図である。

【図 2 0】ping による追放検出手法の一例を示す図である。

【発明を実施するための形態】

【0025】

以下、図面を参照して実施の形態を説明する。

〔1〕一実施形態

〔1-1〕分散ファイルシステムについて

上述したように、クライアント（以下、クライアント A という）が何らかの理由でサーバからの要求に対して反応しなくなった場合、クライアント A に付与された分散ロックは資源排他管理用サブシステムに対して返却されなくなる。その結果、分散ファイルシステム全体（サーバ及び他の全てのクライアント）が、当該分散ロックによって排他されているファイルシステム資源に対して処理が行えない状態に陥る。

20

【0026】

一方、資源排他管理用サブシステムが、クライアント A が所持していた分散ロックを他のクライアント（以下、クライアント B という）に渡すこと等により強制的に処理を進めた場合、クライアント A が復帰して処理を再開することがある。このとき、クライアント A は、新たにロックを付与されたクライアント B と競合を引き起こして、データ破壊を発生させる等の可能性がある。

【0027】

この状況を防ぐため、サーバは、分散ロックを付与されたクライアント A がダウンしたと判断した場合、その時点でサーバからクライアント A に関する全ての情報を無効化する。例えば L u s t r e においては、クライアント A に付与された分散ロックやキャッシュされていた i n o d e の無効化、キャッシュされていたデータがあれば全てのデータのフラッシュ、等を行ない、クライアント A との間の接続を切断する。すなわち、サーバは、クライアントを追放する。

30

【0028】

例えば、クライアント A が動作している装置がダウンしたためにクライアント A がサーバからの要求に反応しなくなった場合、ダウン前にクライアント A が保持していたクライアント側の接続情報は失われてしまう。よって、当該装置の再起動によりクライアントが再起動し、当該クライアントが接続リクエストをサーバへ送信した際には、システム（サーバ）では、新しいクライアント（以下、クライアント C という）が接続リクエストを送信したものとして処理される。このように、装置が、ダウン、再起動、そしてクライアント C のマウントを実行するまでの間に、サーバ側で既にクライアント A の追放が発生していたとしても、クライアント A に関する情報はサーバ上及びクライアント上の双方から削除されている。従って、上述のように、装置がダウンした場合には、分散ファイルシステムにおけるデータ整合性は維持される。

40

【0029】

しかし、サーバが、クライアント A がダウンしたと判断してクライアント追放処理を実行したものの、実際にはクライアント A はダウンしていない場合もある。このような場合としては、例えば、ネットワークエラーによってクライアント A - サーバ間で一時的に通

50

信不可能になっていた場合や、クライアント A 側で CPU の高負荷又はメモリ不足による通信処理不能に陥っていた場合が挙げられる。この場合、クライアント A は、接続情報を保持しているが、当該接続情報は既にサーバ側には存在しない情報である。つまり、クライアント A - サーバ間で接続情報に関して整合が取れていない状態となる。

【 0 0 3 0 】

従って、クライアント A がサーバに存在しない接続情報を保持した状態（サーバと整合が取れていない状態）になった場合、当該接続情報を新しく書き換えるために、何らかの契機で、クライアント A は追放復旧処理を行なうことが好ましい。

そこで、一実施形態に係る分散ファイルシステム 1 は、以下に詳述する処理を行なう。

〔 1 - 2 〕 分散ファイルシステムの構成

以下、図 1 及び図 2 を参照して、一実施形態の一例としての分散ファイルシステム（情報処理システム）1 の構成について説明する。

【 0 0 3 1 】

図 1 は、一実施形態の一例としての分散ファイルシステム 1 の構成例を示す図であり、図 2 は、図 1 に示す F S（File System）クライアント 3 0 に着目した分散ファイルシステム 1 の構成例を示す図である。

図 1 に示すように、分散ファイルシステム 1 は、M G S（Management Server）1 0 - 1、M D S（Meta Data Server）1 0 - 2、及び $n - 2$ 個（ n は 2 以上の整数）の O S S（Object Storage Server）1 0 - 3 ~ 1 0 - n をそなえる。また、分散ファイルシステム 1 は、M G T（Management Target）2 0 - 1、M D T（Meta Data Target）2 0 - 2、及び $n - 2$ 個の O S T（Object Storage Target）2 0 - 3 ~ 2 0 - n をさらにそなえる。さらに、分散ファイルシステム 1 は、 m 個（ m は 0 以上の整数）の F S クライアント 3 0 - 1 ~ 3 0 - m 、及びネットワーク 4 0 をさらにそなえる。

【 0 0 3 2 】

以下、M G S 1 0 - 1、M D S 1 0 - 2、O S S 1 0 - 3 ~ 1 0 - n を区別しない場合には、単にサーバ・サービス（サーバ、情報処理装置）1 0 という。また、M G T 2 0 - 1、M D T 2 0 - 2、O S T 2 0 - 3 ~ 2 0 - n を区別しない場合には、単に論理ボリューム 2 0 という。

分散ファイルシステム 1 としては、L u s t r e - 1 . 8 版を用いたファイルシステムが例として挙げられる。

【 0 0 3 3 】

分散ファイルシステム 1 が想定するシステムは、ネットワーク層を介したサーバ数対 F S クライアント数が n 対 m のシステムである。図 1 に示すように、複数のサーバ 1 0 及び F S クライアント 3 0 は、ネットワーク層（ネットワーク 4 0）の上に位置し、互いに通信し合いながら 1 つの分散ファイルシステム 1 を構成する。

図 1 に示すように、分散ファイルシステム 1 における複数のサーバ 1 0 は、それぞれ 1 つの独自論理ボリューム 2 0 を管理する。

【 0 0 3 4 】

M G T 2 0 - 1 は、分散ファイルシステム 1 の構成情報を保持する論理ボリュームであり、M G S 1 0 - 1 は、M G T 2 0 - 1 を管理するサーバである。

O S T 2 0 - 3 ~ 2 0 - n は、テキストや計算結果等のデータ（ファイル、オブジェクト）を保持する論理ボリュームであり、O S S 1 0 - 3 ~ 1 0 - n は、それぞれ O S T 2 0 - 3 ~ 2 0 - n を管理するサーバである。

【 0 0 3 5 】

M D T 2 0 - 2 は、O S T 2 0 - 3 ~ 2 0 - n が保持するファイルの更新時間やファイルサイズ等のメタデータを保持する論理ボリュームであり、M D S 1 0 - 2 は、M D T 2 0 - 2 を管理するサーバである。

各 F S クライアント 3 0 は、図 2 に示すように、サーバ数に応じたクライアント・サービス（クライアント）をそなえる。具体的には、各 F S クライアント 3 0 は、M G C（Management Client）1 3 0 - 1、M D C（Meta Data Client）1 3 0 - 2、及び $n - 2$ 個

10

20

30

40

50

のOSC (Object Storage Client) 130 - 3 ~ 130 - nをそなえる。

【0036】

以下、MGC 130 - 1、MDC 130 - 2、OSC 130 - 3 ~ 130 - nを区別しない場合には、単にクライアント・サービス(クライアント, 端末装置) 130という。

MGC 130 - 1はMGS 10 - 1との間でのリクエスト処理、MDC 130 - 2はMDS 10 - 2との間でのリクエスト処理、OSC 130 - 3 ~ 130 - nはそれぞれOSS 10 - 3 ~ 10 - nとの間でのリクエスト処理を担っている。

【0037】

つまり、図1に示すサーバ数がn個、FSクライアント数がm個の分散ファイルシステム1においては、サーバ数対クライアント数はn対($n * m$)になる(図2参照)。

ここで、1つのクライアント130に着目すると、クライアント130が対象とする通信相手は、1つのサーバ10のみとなる。以下、サーバ10及びクライアント130の1対1の接続確立方法に着目して説明する。

【0038】

〔1-3〕分散ファイルシステムの説明

ここで、一実施形態に係る分散ファイルシステム1について、簡単に説明する。

上述のように、追放済クライアントにおいて、ユーザ・アプリケーションからのリクエストがサーバへ送信されると、サーバからユーザ・アプリケーションへエラーが返ってしまう。なお、追放済クライアントとは、サーバ側でクライアント追放後に一切リクエストを発行しておらず、よって追放復旧処理が実行されていないクライアントをいう。

【0039】

また、図20に示す手法では、システム規模が拡大し、サーバ数やクライアント数が増大するに従って、定期的にpingリクエストを送信することによる情報処理システムの負荷が膨大になる。

これに対し、一実施形態に係る分散ファイルシステム1は、以下の(i) ~ (iii)の処理を行なうことで、上述した不都合を解消する。

【0040】

(i) サーバ(情報処理装置) 10は、クライアント(端末装置) 130との接続を解除する追放予定時刻(予定時刻)をクライアント130に通知する。

(ii) クライアント130は、サーバ10へリクエスト(要求)を送信する際に、現在時刻がサーバ10から通知された追放予定時刻を経過しているか否かを判断する。

(iii) クライアント130は、現在時刻が追放予定時刻を経過していると判断した場合、サーバ10へリクエストを送信する前に、サーバ10との間で接続を確立するための接続リクエストをサーバ10へ送信する。

【0041】

上記(i) ~ (iii)の処理により、分散ファイルシステム1は、同一の追放予定時刻をサーバ10側及びクライアント130側の双方に持たせ、追放予定時刻を使用してサーバ10側とクライアント130側を同期させることができる。

このように、分散ファイルシステム1は、追放予定時刻を過ぎてサーバ10側で追放済状態のクライアント130に対して、リクエストの送信の際に自身がサーバ10により追放済であることを認識させる(追放復旧処理を行なう契機を与える)ことができる。従って、クライアント130は、リクエストを送信する前に、追放復旧処理により新規接続を確立してから(サーバ10側と整合性のとれた新規の接続情報を作成してから)、リクエストを送信することができる。

【0042】

これにより、クライアント130では、サーバ10では既に削除されたサーバ側接続情報に対応する古いクライアント側接続情報を参照したことによる、重大エラーの発生を抑制でき、ユーザ・アプリケーションに対してエラーが返る確率を下げるができる。

また、定期的なpingリクエストの送信も行なわれないため、図20に示す技術よりも分散ファイルシステム1の負荷(処理負荷)を軽減させることもできる。

10

20

30

40

50

【 0 0 4 3 】

以上のように、分散ファイルシステム 1 によれば、ユーザ・アプリケーションに対してエラーが返される頻度を低減させるとともに、システム 1 の負荷の増大を抑制させながら、サーバ 1 0 とクライアント 1 3 0 との間の接続を再確立することができる。

〔 1 - 4 〕 ハードウェア構成

次に、図 3 を参照して、分散ファイルシステム 1 のハードウェア構成について説明する。図 3 は、図 1 に示すサーバ 1 0 及び F S クライアント 3 0 それぞれのハードウェア構成例を示す図である。

【 0 0 4 4 】

サーバ 1 0 及び F S クライアント 3 0 は、図 3 に示すように、それぞれ C P U 1 0 a 及び 3 0 a、メモリ 1 0 b 及び 3 0 b、記憶部 1 0 c 及び 3 0 c、ネットワークインタフェース 1 0 d 及び 3 0 d、並びに入出力部 1 0 e 及び 3 0 e をそなえる。また、サーバ 1 0 及び F S クライアント 3 0 は、図 3 に示すように、記録媒体 1 0 f 及び 3 0 f、並びに読取部 1 0 g 及び 3 0 g をさらにそなえる。

【 0 0 4 5 】

C P U 1 0 a 及び 3 0 a は、それぞれ、図 3 における対応する各ブロック 1 0 b ~ 1 0 g 及び 3 0 b ~ 3 0 g と接続され、種々の制御や演算を行なう処理装置（プロセッサ）である。C P U 1 0 a 及び 3 0 a は、メモリ 1 0 b 及び 3 0 b、記録媒体 1 0 f 及び 3 0 f、又は図示しない R O M（Read Only Memory）等に格納されたプログラムを実行することにより、サーバ 1 0 又は F S クライアント 3 0 における種々の機能を実現する。

【 0 0 4 6 】

メモリ 1 0 b 及び 3 0 b は、種々のデータやプログラムを一時的に格納する記憶装置である。C P U 1 0 a 及び 3 0 a は、プログラムを実行する際に、メモリ 1 0 b 及び 3 0 b にデータやプログラムを格納し展開する。なお、メモリ 1 0 b 及び 3 0 b としては、例えば R A M（Random Access Memory）等の揮発性メモリが挙げられる。

記憶部 1 0 c 及び 3 0 c は、種々のデータやプログラム等を格納するハードウェアである。記憶部 1 0 c 及び 3 0 c としては、例えば H D D（Hard Disk Drive）等の磁気ディスク装置、S S D（Solid State Drive）等の半導体ドライブ装置、フラッシュメモリ等の不揮発性メモリ等の各種デバイスが挙げられる。

【 0 0 4 7 】

ネットワークインタフェース部 1 0 d 及び 3 0 d は、有線又は無線による、サーバ 1 0 - ネットワーク 4 0 間、又は、F S クライアント 3 0 - ネットワーク 4 0 間の接続及び通信の制御を行なう。ネットワークインタフェース部 1 0 d 及び 3 0 d としては、T C P（Transmission Control Protocol）/ I P をサポートする L A N（Local Area Network）カード等のネットワークコントローラが例として挙げられる。また、ネットワークインタフェース部 1 0 d 及び 3 0 d としては、インフィニバンド（InfiniBand、登録商標）等の H C A（Host Channel Adapter）や、ファイバチャネル（Fibre Channel）コントローラ等も例として挙げられる。

【 0 0 4 8 】

入出力部 1 0 e 及び 3 0 e は、例えばマウスやキーボード等の入力装置、及びディスプレイやプリンタ等の出力装置の少なくとも一方を含むものである。入出力部 1 0 e 及び 3 0 e は、入力装置によりサーバ 1 0 及び F S クライアント 3 0 のオペレータ（管理者）やユーザ（使用者）の操作等による動作命令を受け付ける一方、サーバ 1 0 及び F S クライアント 3 0 による監視結果等の処理結果を出力装置に表示（出力）する。

【 0 0 4 9 】

記録媒体 1 0 f 及び 3 0 f は、フラッシュメモリや R O M 等の記憶装置であり、種々のデータやプログラムを記録する。読取部 1 0 g 及び 3 0 g は、光ディスクや U S B（Universal Serial Bus）メモリ等のコンピュータ読取可能な記録媒体 1 0 h 及び 3 0 h に記録されたデータやプログラムを読み出す装置である。

記録媒体 1 0 f 及び 1 0 h の少なくとも一方には、本実施形態に係るサーバ 1 0 の機能

10

20

30

40

50

を実現する制御プログラムが格納されてもよく、記録媒体 30 f 及び 30 h の少なくとも一方には、F S クライアント 30 の機能を実現する制御プログラムが格納されてもよい。例えば、C P U 10 a 及び 30 a は、それぞれ、記録媒体 10 f 及び 30 f から読み出した制御プログラム、又は、読取部 10 g 及び 30 g を介して記録媒体 10 h 及び 30 h から読み出した制御プログラムを、メモリ 10 b 及び 30 b 等の記憶装置に展開して実行する。これにより、サーバ 10 としてのコンピュータ及び F S クライアント 30 としてのコンピュータは、C P U 10 a 及び 30 a により、本実施形態に係るサーバ 10 及び F S クライアント 30 の機能を実現する。

【0050】

なお、上述した各ブロック 10 a ~ 10 g 間、各ブロック 30 a ~ 30 g 間は、それぞれバスで相互に通信可能に接続される。

10

また、分散ファイルシステム 1 の上述したハードウェア構成は例示である。従って、個々のストレージシステム 1 内、サーバ 10 内、又は F S クライアント 30 内でのハードウェアの増減や分割、任意の組み合わせでの統合等は、適宜行なわれてもよい。例えば、図 3 に示すサーバ 10 のハードウェアは、1 以上のサーバ 10 で共用されてもよく、図 3 に示す F S クライアント 30 のハードウェアは、1 以上の F S クライアント 30 で共用されてもよい。

【0051】

さらに、図 1 に示す論理ボリューム 20 は、複数の物理ボリュームを搭載する図示しないストレージ装置の記憶領域や記憶部 10 c の記憶領域等により実現される。

20

〔1-5〕分散ファイルシステムの詳細な構成

〔1-5-1〕サーバの構成

次に、図 4 を参照して、一実施形態の一例としてのサーバ 10 の構成について説明する。図 4 は、図 1 に示すサーバ 10 の機能構成例を示す図である。

【0052】

サーバ 10 は、上述のように、M G S 10 - 1、M D S 10 - 2、又は O S S 10 - 3 ~ 10 - n としての機能を有し、接続を確立したクライアント 130 に対してサービスを提供する。

また、一実施形態に係るサーバ 10 は、保持部 11、受信処理部 12、リクエスト処理部 13、及び追放処理部 14 を有する。

30

【0053】

保持部 11 は、クライアント 130 ごとに、クライアント 130 に関する接続情報 11 a を保持するものであり、例えば、メモリ 10 b 又は記憶部 10 c 等により実現される。

接続情報 11 a には、接続管理情報及び資源排他管理情報が含まれ得る。接続管理情報は、クライアント 130 との接続（コネクション）に関する情報であり、資源排他管理情報は、排他制御に用いられる構造体（例えばリソースのロック範囲）及びそれに紐付けられたリソースに関する情報である。

【0054】

また、接続情報 11 a は、クライアント 130 を追放する（クライアント 130 との接続を解除する）追放予定時刻（予定時刻）11 b を含む。

40

図 5 は、一実施形態の一例としてのサーバ 10 及びクライアント 130 が保持する追放予定時刻を説明する図である。

例えば、図 5 に示すように、サーバ 10（図 5 の説明では、便宜上サーバ A という）が、3 つのクライアント 130（図 5 の説明では、便宜上クライアント A ~ C という）の各々との間で接続を確立している場合を説明する。この場合、サーバ A は、保持部 11 に、サーバ A と接続を確立している全てのクライアント 130 の接続情報 11 a、つまり、クライアント A の接続情報 11 a - 1、クライアント B の接続情報 11 a - 2、及びクライアント C の接続情報 11 a - 3 を保持する。また、サーバ A が保持する接続情報 11 a - 1 ~ 11 a - 3 には、それぞれ、クライアント A ~ C の追放予定時刻 11 b - 1 ~ 11 b - 3 が含まれる。

50

【 0 0 5 5 】

受信処理部 1 2 は、クライアント 1 3 0 からの接続リクエスト、書込 / 読出等の各種リクエスト等の情報を受信し、接続リクエストや各種リクエストが、サーバ 1 0 が保持する接続情報 1 1 a に従っているか否かを判断して、判断結果に応じた所定の処理を行なう。

例えば、受信処理部 1 2 は、接続リクエストを受信すると、接続リクエストの送信元のクライアント 1 3 0 に関する接続情報 1 1 a が保持部 1 1 に保持されているか否かを判断する。保持されていない場合、受信処理部 1 2 は、当該クライアント 1 3 0 に関する接続情報 1 1 a を新規に作成するため、接続リクエストをリクエスト処理部 1 3 に渡す。

【 0 0 5 6 】

一方、当該クライアント 1 3 0 に関する接続情報 1 1 a が保持部 1 1 に保持されている場合、受信処理部 1 2 は、現在時刻が追放予定時刻 1 1 b を経過しているか否かを判断する。現在時刻が追放予定時刻 1 1 b を経過していない場合、リクエスト処理部 1 3 に当該クライアント 1 3 0 へのエラーを返させる。また、現在時刻が追放予定時刻 1 1 b を経過している場合、受信処理部 1 2 は、当該クライアント 1 3 0 に関する既存の接続情報 1 1 a を削除するため、追放処理部 1 4 に所定のシグナルを送信する。シグナル送信後、受信処理部 1 2 は、当該クライアント 1 3 0 に関する接続情報 1 1 a を新規に作成するため、接続リクエストをリクエスト処理部 1 3 に渡す。

【 0 0 5 7 】

また、受信処理部 1 2 は、クライアント 1 3 0 から、接続リクエスト以外の書込 / 読出等の各種リクエストを受信すると、当該リクエストの送信元のクライアント 1 3 0 に関する接続情報 1 1 a が保持部 1 1 に保持されているか否かを判断する。保持されていない場合、受信処理部 1 2 は、リクエスト処理部 1 3 に当該クライアント 1 3 0 へのエラーを返させる。

【 0 0 5 8 】

一方、当該クライアント 1 3 0 に関する接続情報 1 1 a が保持部 1 1 に保持されている場合、受信処理部 1 2 は、当該クライアント 1 3 0 からのリクエストをリクエスト処理部 1 3 に渡す。

リクエスト処理部 1 3 は、接続情報 1 1 a 及び追放予定時刻 1 1 b の管理、並びに受信処理部 1 2 が受信した接続リクエスト又は各種リクエストに応じた処理、リプライの作成及び送信等を行なう。

【 0 0 5 9 】

例えば、リクエスト処理部（接続管理部）1 3 は、受信処理部 1 2 から接続リクエスト又は各種リクエストを渡されると、追放予定時刻 1 1 b を取得する。そして、リクエスト処理部 1 3 は、接続リクエスト又は各種リクエストの送信元のクライアント 1 3 0 に関する接続情報 1 1 a と対応付けて管理する。

具体的には、リクエスト処理部 1 3 は、受信処理部 1 2 から接続リクエストを渡されると、当該クライアント 1 3 0 に関する新規の接続情報 1 1 a を生成し、生成した接続情報 1 1 a を、当該クライアント 1 3 0 に対応付けて保持部 1 1 に保持させる。そして、リクエスト処理部 1 3 は、接続リクエストに応じた所定の処理を実行し、追放予定時刻 1 1 b を計算して、生成した接続情報 1 1 a に記録する。

【 0 0 6 0 】

また、リクエスト処理部 1 3 は、受信処理部 1 2 から各種リクエストを渡されると、当該リクエストに応じた所定の処理（例えば書込処理 / 読出処理等）を実行し、追放予定時刻 1 1 b を更新して、当該クライアント 1 3 0 に関する接続情報 1 1 a に記録する。

そして、リクエスト処理部 1 3 は、データ部に各種情報を記録した、接続リクエスト又は各種リクエストへのリプライを生成し、生成したリプライを接続リクエスト又は各種リクエストの送信元のクライアント 1 3 0 へ送信する。

【 0 0 6 1 】

ここで、追放予定時刻 1 1 b は、クライアント 1 3 0 から接続リクエスト又は各種リクエストを受けると、リクエスト処理部 1 3 により、現在時刻に所定時間（例えば 2 5 秒）

10

20

30

40

50

が加算されることにより取得（計算・更新）される。

また、リクエスト処理部（通知部）13は、取得した追放予定時刻11bを、接続リクエスト又は各種リクエストの送信元のクライアント130へ送信する。

【0062】

図6は、図4に示すリクエスト処理部13による追放予定時刻11bの通知手法の一例を示す図である。

例えば、サーバ10（リクエスト処理部13）は、図6に示すように、クライアント130から受信した接続リクエスト又は各種リクエストに対するリプライ（例えばヘッダ部又はデータ部）に、追放予定時刻11bを設定して（含めて）送信することができる。つまり、リクエスト処理部13による追放予定時刻11bの取得処理は、クライアント130に対してリプライ（応答）を送信（作成）する際に行なわれてよい。追放予定時刻11bがリプライを送信（作成）する際に取得される場合、現在時刻は、リプライの送信（作成）時の時刻となる。

【0063】

なお、図6に示すように、接続リクエスト又は各種リクエストには、ヘッダ部に送信先及び送信元を一意に定めるために使用される送信先ID及び送信元IDが含まれる。また、サーバ10は、これらの情報（例えば接続リクエスト又は各種リクエスト中の送信元ID）に基づいて、処理対象となる接続情報11aを一意に求めることができる。

以上のように、リクエスト処理部13は、クライアント130から初めて接続リクエストを受信すると、当該クライアント130に関する接続情報11aを生成して保持部11に保持する。また、リクエスト処理部13は、取得（算出）した追放予定時刻11bを接続情報11aに記録するとともに、当該クライアント130へ通知する。なお、「クライアント130から初めて接続リクエストを受信」とは、追放済且つ追放復旧処理が行なわれていないクライアント130等、サーバ10が対応する接続情報11aを保持していないクライアント130から接続リクエストを受信した場合も含む。

【0064】

また、リクエスト処理部13は、クライアント130から各種リクエストを受信するたびに、取得（更新）した追放予定時刻11bを接続情報11aに記録するとともに、当該クライアント130へ通知する。

なお、リクエスト処理部13は、受信処理部12からエラーの送信を指示されると、接続リクエスト又は各種リクエストの送信元のクライアント130へ、エラーを送信する。エラーが送信される原因としては、上述のように、クライアント130に関する接続情報11aが保持部11に保持されている状態で、当該クライアント130から接続リクエストを受けた場合が挙げられる。また、エラーが送信される他の原因としては、保持部11に接続情報11aが保持されていないクライアント130から、接続リクエスト以外の各種リクエストを受けた場合等も挙げられる。

【0065】

追放処理部（接続解除部）14は、現在時刻が追放予定時刻11bを経過しているクライアント130を追放する（クライアント130との間で確立された接続を解除する）クライアント追放処理（接続解除処理）を行なう。

具体的には、追放処理部14は、所定時間（例えば25秒）ごとに、保持部11に保持された1以上の接続情報11aを順に参照し、現在時刻が追放予定時刻11bを経過しているか否かを判断する。そして、追放処理部14は、現在時刻が追放予定時刻11bを経過していると判断した場合、対応する接続情報11aを無効化し、クライアント130との間の接続を切り離す（切断する）。

【0066】

なお、リクエスト処理部13は、例えば、保持部11から当該クライアント130に関する接続管理情報及び資源排他情報等の接続情報11aを削除することで、接続情報11aを無効化することができる。

また、追放処理部14は、受信処理部12から所定のシグナルを受信した場合にも、保

10

20

30

40

50

持部 11 に保持された 1 以上の接続情報 11a についてクライアント追放処理を行なう。

【0067】

すなわち、追放処理部 14 は、所定のタイミングとして、前回クライアント追放処理を行なってから所定時間が経過したとき、及び所定のシグナルを受信したとき、のいずれか早いタイミングで、クライアント追放処理を行なう。

なお、保持部 11 に保持された 1 以上の接続情報 11a は、例えば双方向リストにより管理され、追放処理部 14 は、クライアント追放処理において、双方向リストを順に辿ることで、判断対象の接続情報 11a を選択することができる。

【0068】

以上のように、追放処理部 14 は、定期的にクライアント 130 の接続情報 11a を監視し、追放予定時刻 11b を経過したクライアント 130 の接続情報 11a がある場合には当該クライアント 130 を追放する。

従来、クライアント - サーバ間で一度接続が確立されると、クライアントが明示的に接続の切断を指示されるまで、サーバにはクライアントの接続情報が保持されていることが期待されていた。しかし、上述した追放処理部 14 によれば、追放予定時刻 11b を経過したクライアント 130 の接続情報 11a を積極的に破棄することが可能となり、サーバ 10 のメモリ 10b 等のメモリ使用量を削減することができる。また、サーバ 10 上のクライアント 130 の接続情報量が削減されるため、接続情報 11a の検索速度を向上させることができる。

【0069】

また、リクエスト処理部 13 は、設定する追放予定時刻 11b (差分間隔) を、例えば図 20 に示す例の ping リクエストの送信間隔と同程度とすることができる。これにより、分散ファイルシステム 1 は、ping 方式と同程度のレベルで、ユーザ・アプリケーションに重大なエラーが返る可能性を確実に排除することができる。

なお、上述したサーバ 10 において、受信処理部 12 及びリクエスト処理部 13 の機能は、リクエストハンドラを実行するスレッドに持たせることができる。また、追放処理部 14 の機能は、クライアント 130 の追放を検知・実行するスレッド (クライアント追放用スレッド) に持たせることができる。

【0070】

サーバ 10 は、リクエストハンドラを実行するスレッドを複数 (例えばクライアント 130 の数) 実行することができ、また、1 つのクライアント追放用スレッドを実行することができる。例えば、クライアント追放用スレッドは、リクエストハンドラを実行する各スレッドから所定の信号を入力される都度、起動し、クライアント追放処理を実行する。

なお、受信処理部 12 及びリクエスト処理部 13 は、例えば、ネットワークインタフェース部 10d と、メモリ 10b に展開された制御プログラムを実行する CPU 10a とが協働することにより実現される。また、追放処理部 14 は、例えば、メモリ 10b に展開された制御プログラムを実行する CPU 10a により実現される。

【0071】

〔1-5-2〕クライアントの構成

次に、図 7 を参照して、一実施形態の一例としてのクライアント 130 の構成について説明する。図 7 は、図 2 に示すクライアント 130 の機能構成例を示す図である。

クライアント 130 は、上述のように、各 FS クライアント 30 内にサーバ 10 の数と同数そなえられ、対応するサーバ 10 と当該サーバ 10 との間で確立された接続を用いて 1 対 1 の通信を行なう。

【0072】

また、一実施形態に係るクライアント 130 は、保持部 31、受信処理部 32、時刻管理部 33、及び送信処理部 34 を有する。

保持部 31 は、自クライアント 130 に対応するサーバ 10 に関する接続情報 31a を保持するものであり、例えば、メモリ 30b 又は記憶部 30c 等により実現される。

なお、接続情報 31a は、自クライアント 130 が追放済クライアントでない場合、接

10

20

30

40

50

続相手のサーバ10が保持する自クライアント130に関する接続情報11aに対応するものであり、当該接続情報11aと同様の情報を含むことができる。

【0073】

また、接続情報31aは、サーバ10から通知された追放予定時刻（予定時刻）31bを含む。なお、追放予定時刻31bは、対応するサーバ10が管理する追放予定時刻11bと同一の時刻であるが、便宜上、クライアント130が保持するものを追放予定時刻31bという。

例えば、図5に示すように、クライアントA～Cは、それぞれ、保持部31-1～31-3に、サーバAの接続情報31a-1～31a-3を保持する。また、接続情報31a-1～31a-3には、それぞれ、各クライアントA～Cがサーバ10から通知された追放予定時刻31b-1～31b-3が含まれる。なお、追放予定時刻31b-1～31b-3は、それぞれ、サーバ10が保持部11に保持する追放予定時刻11b-1～11b-3と同時刻である。

10

【0074】

受信処理部32は、サーバ10から送信されたリクエスト（例えば分散ロック要求）や、リクエストに対するリプライ、エラー等の情報を受信する。

また、受信処理部32は、接続リクエスト又は各種リクエストを送信したサーバ10から、リプライを受信すると、当該リプライに含まれる追放予定時刻31bを取得して、管理部33に渡す。

【0075】

20

管理部33は、接続情報31a及び追放予定時刻31bの管理、並びに各種リクエストの送信の際に現在時刻と追放予定時刻31bとの比較を行なう。

例えば、管理部（時刻管理部）33は、受信処理部32がサーバ10から受信した追放予定時刻31bをサーバ10に関する接続情報31aに対応付けて管理する。

また、管理部（判断部）33は、送信処理部34がサーバ10へ書込・読出等の各種リクエストを送信する際に、現在時刻がサーバ10から通知され管理部33が管理する追放予定時刻31bを経過しているか否かを判断する。

【0076】

現在時刻が追放予定時刻31bを経過していると判断した場合、管理部（接続情報管理部）33は、サーバ10に関する接続情報31aを無効化（例えば保持部31から削除）する。そして、管理部33は、送信処理部34にサーバ10へ接続リクエストを送信させ、サーバ10との間で接続を確立させる。つまり、管理部33は、受信処理部32がサーバ10から接続リクエストのリプライを受信すると、当該リプライの内容に基づき接続情報31aを作成し、保持部31に保持させる。そして、管理部33は、送信処理部34に、作成した接続情報31aに基づいてサーバ10へ各種リクエストを送信させる。

30

【0077】

一方、現在時刻が追放予定時刻31bを経過していないと判断した場合、管理部33は、送信処理部34に、既に保持部31に保持されている接続情報31aに基づいてサーバ10へ各種リクエストを送信させる。

以上のように、管理部33は、自クライアント130がサーバ10側で追放済であると判断すると、送信処理部34に、サーバ10へ各種リクエストを送信させる前に、接続リクエストをサーバ10へ送信させる。

40

【0078】

送信処理部（送信部）34は、管理部33からの指示に応じて、クライアント130上のユーザ・アプリケーションやシステムが生成した接続リクエスト、各種リクエスト等の情報をサーバ10へ送信する。例えば、送信処理部34は、図6に示すように、送信先IDとしてのサーバ10のID及び送信元IDとしてのクライアント130のID、その他接続情報31a等を含む接続リクエスト／各種リクエストを生成して送信する。

【0079】

なお、クライアント130では、接続リクエスト及び各種リクエストの発行元が、ユー

50

ザ・アプリケーションからシステムスレッドまで多様である。例えば、接続リクエストは、クライアント130のシステム等から発行され、各種リクエストは、クライアント130上で実行されるユーザ・アプリケーションによるシステムコール(s y s c a l l)等である。

【0080】

以上のように、クライアント130は、リクエスト送信のたびに、サーバ10から通知された追放予定時刻31bを参照し、現在時刻が追放予定時刻31bを超過していた場合、これまでの接続情報31aを使用せず、再接続処理が完了してからリクエストを送信する。

図20に示す例では、クライアントのpingリクエスト送信先はシステムを構成するサーバ数に比例して増加するため、サーバ数が増えるほど、各クライアントでpingリクエスト送信処理のためのCPU負荷が増大する。これに対し、クライアント130によれば、サーバ10側の追放検知のためにクライアント130側からpingリクエストを送信せずに済むため、ping方式のようにサーバ台数に比例してクライアントのCPU負荷が増加することを抑制できる。

【0081】

また、図20に示す例では、pingを受けたサーバは送信元のクライアントにpingリクエストに対する返信を行なう。サーバが返信するpingリクエストの量は、クライアント数に比例して増加するため、クライアント数が増えるほど、各サーバでpingリクエスト返信処理のためのCPU負荷が増大する。これに対し、クライアント130によれば、サーバ10側はクライアント130から受信したpingリクエストに対する返信を行わずに済むため、ping方式のようにクライアント台数に比例してサーバのCPU負荷が増加することを抑制できる。

【0082】

さらに、図20に示す例のping方式では、pingリクエストの送信及び返信にネットワーク資源が使用される。上述のように、pingリクエストによるネットワークの負荷は、ping負荷*クライアント数*サーバ数になる。よって、システムが大規模化すると、pingリクエストによるネットワーク負荷は増大し、ネットワーク帯域を圧迫することになる。これに対し、クライアント130によれば、クライアント130の追放復旧処理にpingリクエストを送信せずに済むため、ping方式のようにシステム規模(サーバ及びクライアント台数)に比例してシステムのネットワーク負荷が増大することを抑制できる。

【0083】

なお、受信処理部32及び送信処理部34は、例えば、ネットワークインタフェース部30dと、メモリ30bに展開された制御プログラムを実行するCPU30aとが協働することにより実現される。また、管理部33は、例えば、メモリ30bに展開された制御プログラムを実行するCPU30aにより実現される。

〔1-5-3〕サーバ及びクライアント間の通信

次に、図8を参照して、図1に示すサーバ10及びクライアント130間の通信について説明する。図8は、図2に示すサーバ10及びクライアント130間の通信の一例を説明するシーケンス図である。

【0084】

図8に示すように、サーバ10(リクエスト処理部13)は、各クライアント130から最後に受信したリクエストの時刻から、追放予定時刻11bを計算(又は更新)し、クライアント130に追放予定時刻11b(31b)を通知する(処理T1, T2)。このとき、サーバ10は、追放予定時刻11bとして例えば受信時刻又は返信時刻(図8の例では返信時刻)から25秒後を計算して、保持部11に記録する。

【0085】

なお、クライアント130は、サーバ10から通知された追放予定時刻31bを保持部31に記録する。クライアント130(送信処理部34)は、接続リクエスト以外の各種

10

20

30

40

50

リクエストを送信する際に、追放予定時刻 3 1 bを確認し、現在時刻が追放予定時刻 3 1 bを経過していない場合、各種リクエストを送信する。

ここで、サーバ 1 0（追放処理部 1 4）は、追放予定時刻 1 1 bがくると、クライアント 1 3 0をサーバ 1 0から追放する（処理 T 3）。

【 0 0 8 6 】

クライアント 1 3 0（送信処理部 3 4）は、各種リクエスト、例えばユーザ・アプリケーションによるシステムコールを送信する際に、現在時刻が追放予定時刻 3 1 bを経過していた場合、当該リクエストを送信する前に、接続リクエストを送信する（処理 T 4）。サーバ 1 0は、接続リクエストに応じて再接続処理を行なって接続状態を確立し、新たに計算した追放予定時刻 3 1 bをリプライに含めてクライアント 1 3 0へ送信する（処理 T 5）。

10

【 0 0 8 7 】

クライアント 1 3 0は、再接続が確立されてから、サーバ 1 0へリクエスト（システムコール）を送信する（処理 T 6）。サーバ 1 0（リクエスト処理部 1 3）は、クライアント 1 3 0からのリクエストに対して処理を行ない、リプライを返す（処理 T 7）。

このように、サーバ 1 0は、クライアント 1 3 0からのリクエストを受信した際に取得した追放予定時刻 1 1 bに基づいて、クライアント 1 3 0接続情報 3 1 aを追放（無効化）するか否かを決定する。また、クライアント 1 3 0は、サーバ 1 0から通知された追放予定時刻 3 1 bに基づいて、通常のリクエスト（各種リクエスト）の送信前に、接続リクエストを送信するか否かを決定する。

20

【 0 0 8 8 】

ところで、既述のように、サーバ側のクライアント追放処理と、クライアント側の追放復旧処理とは、完全に非同期的に行なわれるため、追放済クライアントが長期間に亘って存在することもしばしばある。

追放済クライアントでは、自身がサーバ側で追放されていることを認識していない状態で、起動されたユーザ・アプリケーションから追放済クライアントへのリクエストが発行される場合がある。追放済クライアントは、発行されたリクエストがサーバへのアクセスを伴うものであると、リクエストをサーバへ送信するが、サーバは、追放済クライアントからのリクエストに対して重大エラーを返す。このように、ユーザ・アプリケーションからのリクエストが、追放済クライアントへ上述した追放復旧契機を与えしまい、ユーザ・アプリケーションに対してエラーが返ってしまうことがしばしば発生する。

30

【 0 0 8 9 】

以上のように、クライアント追放処理は、システム中のデータ破壊やシステム自体の破壊の可能性を排除しつつ、システム全体がハング状態に陥る可能性を排除する機能である。つまり、クライアント追放処理が実行されたためにユーザ・アプリケーションに対してエラーが返ってしまう可能性を、完全に排除することは難しい。

しかしながら、追放済状態になってから数分～数時間経過しているような場合でも追放済クライアントが追放復旧処理を実行せず、ユーザ・アプリケーションが発行したリクエストにサーバからエラーが返ってしまうという状況はシステムの安定性からみて好ましくない。

40

【 0 0 9 0 】

これに対して、一実施形態に係る分散ファイルシステム 1では、クライアント 1 3 0が追放予定時刻 3 1 bを持つ。これにより、クライアント 1 3 0は、各種リクエストを送信する際に、自クライアント 1 3 0が追放済か否かを判断し、追放済である場合、再接続を行なって新規の接続情報 3 1 aを獲得してから当該リクエストを送信する。よって、上述したように、長期間に亘って追放済状態であった古い接続情報 3 1 aを参照してクライアント 1 3 0が各種リクエストを送信した結果、重大なエラーがユーザ・アプリケーションに返されるという状況の発生を抑止することができる。

【 0 0 9 1 】

〔 1 - 6 〕動作例

50

次に、図 9 ～ 図 14 を参照して、上述の如く構成された一実施形態の一例としての分散ファイルシステム 1 における動作例を説明する。

〔 1 - 6 - 1 〕サーバ側の動作例

はじめに、サーバ 10 による動作例を説明する。図 9 及び図 10 は、図 2 に示すサーバ 10 によるリクエスト受信処理及びクライアント追放処理の一例をそれぞれ説明するフローチャートである。

【 0092 】

〔 1 - 6 - 1 - 1 〕リクエスト受信処理

サーバ 10 によるリクエスト受信処理では、リクエストハンドラを実行するスレッド（受信処理部 12）により、サーバ 10 に到着したリクエストが受け取られ、リクエストごとに固有の処理が行なわれてから、処理結果が送信元クライアント 130 へ返信される。

10

具体的には、図 9 に示すように、サーバ 10（受信処理部 12）により、クライアント 130 からリクエストが到着するまで待ち合わせが行なわれる（ステップ S1）。リクエストが受信されると、受信処理部 12 により、受信したリクエストの内容が確認され、例えば、リクエストが接続リクエストであるか否かが判断される（ステップ S2）。

【 0093 】

リクエストが接続リクエストであると判断された場合（ステップ S2 の Yes ルート）、受信処理部 12 により、リクエストに記録されているクライアント 130 側の接続情報 31a からサーバ 10 上の接続情報 11a が検索される。そして、受信処理部 12 により、サーバ 10 上に送信元クライアント 130 に関する接続情報 11a が存在するか否かが判断される（ステップ S3）。

20

【 0094 】

接続情報 11a が存在すると判断された場合（ステップ S3 の Yes ルート）、既に送信元クライアント 130 と接続が確立しているため、受信処理部 12 により、現在時刻が追放予定時刻 11b を経過しているか否かが判断される（ステップ S4）。すなわち、既に接続情報 11a が存在し、クライアント 130 が追放予定時刻 11b 通りに接続要求を送っていたとしても、サーバ 10 での追放処理のタイミングが遅れてクライアント 130 の接続情報 11a が残ってしまう等の状況が想定される。なお、サーバ 10 での追放処理のタイミングが遅れる場合としては、サーバ 10 で処理負荷が高くて処理が滞った場合や、通信で想定以上の時間がかかってしまった場合等が挙げられる。

30

【 0095 】

そこで、このような遅延状況を考慮し、サーバ 10 は、接続確立状態で接続リクエストを受信した場合、上記ステップ S4 の処理を行なう。現在時刻が追放予定時刻 11b を経過していると判断された場合（ステップ S4 の Yes ルート）、受信処理部 12 により、クライアント追放処理を実行させるためにクライアント追放用スレッド（追放処理部 14）にシグナルが送信される（ステップ S5）。

【 0096 】

次いで、追放処理部 14 による後述するクライアント追放処理により、新規に接続情報 11a が生成され、保持部 11 に記録される（ステップ S6）。そして、リクエスト処理部 13 により、各種リクエストに固有の処理が行なわれ、処理結果がリクエストに対する返信メッセージ（リプライ）に設定される（ステップ S7）。

40

また、リクエスト処理部 13 により、現在時刻に x 秒（例えば 25 秒）を加算した追放予定時刻 11b が作成され、接続情報 11a に記録される（ステップ S8）。そして、リクエスト処理部 13 により、追放予定時刻 11b（31b）がステップ S7 で生成された返信メッセージに設定され、クライアント 130 へ送信され（ステップ S9）、処理がステップ S1 に移行する。

【 0097 】

一方、ステップ S3 において、受信処理部 12 により接続情報 11a が存在しないと判断された場合（ステップ S3 の No ルート）、送信元クライアント 130 と接続が確立されていない（又は追放済である）ため、処理がステップ S6 に移行する。

50

また、ステップS 4において、現在時刻が追放予定時刻1 1 b以前であると判断された場合（ステップS 4のN o ルート）、処理がステップS 1 1に移行する。すなわち、この場合、接続が確立されているのにクライアント1 3 0から接続リクエストが送信されてきたという状況であるため、リクエスト処理部1 3により、接続が既に確立されていることを表すエラーが送信元クライアント1 3 0に返信される。そして、処理がステップS 1 1に移行する。

【0 0 9 8】

さらに、ステップS 2において、リクエストが接続リクエスト以外の各種リクエストであると判断された場合（ステップS 2のN o ルート）、受信処理部1 2により、リクエストに記録されているクライアント1 3 0側の接続情報3 1 aからサーバ1 0上の接続情報1 1 aが検索される。そして、受信処理部1 2により、サーバ1 0上に送信元クライアント1 3 0に関する接続情報1 1 aが存在するか否かが判断される（ステップS 1 0）。 10

【0 0 9 9】

接続情報1 1 aが存在しないと判断された場合（ステップS 1 0のN o ルート）、接続が確立されていない状態で各種リクエストが受信されたため、リクエスト処理部1 3により、クライアント1 3 0へエラーが返される（ステップS 1 1）。そして、処理がステップS 1 1に移行する。

一方、接続情報1 1 aが存在すると判断された場合（ステップS 1 0のY e s ルート）、リクエスト固有の処理を進めるために処理がステップS 7に移行する。 20

【0 1 0 0】

なお、この場合にも、ステップS 4について上述した説明と同様に、現在時刻が追放予定時刻1 1 bを経過している可能性がある。しかし、接続リクエスト以外の各種リクエストに関しては、ステップS 1 0の判断の段階では接続情報1 1 aが参照できている。従って、後述するクライアント追放用スレッドにおいて、リクエスト受信処理中に接続情報1 1 aが削除されないように排他関係が適切にとられていれば、リクエスト受信処理完了時に追放予定時刻1 1 bが上書き更新される。これにより、クライアント1 3 0が追放されなくなるため、上述したステップS 4と同様の問題は発生しない。例えば、リクエストハンドラを実行するスレッドと、クライアント追放用スレッドとの間で、同時に同じ接続情報1 1 aの処理を行なわないように、スピンロック等をつけること等により、適切な排他関係とすることができる。 30

【0 1 0 1】

〔1 - 6 - 1 - 2〕クライアント追放処理

クライアント追放用スレッド（追放処理部1 4）は、シグナル若しくはx秒（例えば2 5秒）間隔で動作開始し、接続情報1 1 aの追放予定時刻1 1 bを確認して、現在時刻を超過しているものがあればクライアント追放処理を実行する。

具体的には、図1 0に示すように、追放処理部1 4により、シグナルが受信される若しくは2 5秒が経過するまで待機され（ステップS 2 1, S 2 2, S 2 2のN o ルート）。シグナルが受信される若しくは2 5秒が経過すると（ステップS 2 2のY e s ルート）、追放処理部1 4により、未判定の接続情報1 1 aの追放予定時刻1 1 bが取得され（ステップS 2 3）、現在時刻が追放予定時刻1 1 bを経過しているか否かが判断される（ステップS 2 4）。 40

【0 1 0 2】

現在時刻が追放予定時刻1 1 bを経過していると判断された場合（ステップS 2 4のY e s ルート）、追放処理部1 4により、接続情報1 1 aが削除される（ステップS 2 5）。そして、追放処理部1 4により、保持部1 1が保持する全ての接続情報1 1 aに対して判定が行なわれたか否かが判断される（ステップS 2 6）。

全ての接続情報1 1 aに対して判定が行なわれたと判断された場合（ステップS 2 6のY e s ルート）、処理がステップS 2 1に移行する。一方、全ての接続情報1 1 aに対して判定が行なわれていないと判断された場合（ステップS 2 6のN o ルート）、処理がステップS 2 3に移行する。 50

【 0 1 0 3 】

また、ステップ S 2 4 において、現在時刻が追放予定時刻 1 1 b を経過していないと判断された場合（ステップ S 2 4 の N o ルート）、処理がステップ S 2 6 に移行する。

ところで、クライアント追放処理では、接続情報 1 1 a について、ステップ S 2 4 の確認が行なわれた直後に現在時刻が追放予定時刻 1 1 b を経過してしまう場合が考えられるが、以下の（ a ）及び（ b ）により問題とはならない。

【 0 1 0 4 】

（ a ）この接続情報 1 1 a を対象にした接続リクエストが到着すると、図 9 のステップ S 3 の判断により古い接続情報 1 1 a が使用されず、新規接続情報 1 1 a が作成される（ステップ S 3 の N o ルート， S 6 ）。また、ステップ S 5 で送信されるシグナルによって、クライアント追放用スレッドが現在のループ処理を完了してもすぐに起床してループ処理を再開することになるため（図 1 0 のステップ S 2 2 の Y e s ルート）、古い接続情報 1 1 a は速やかに削除される。

10

【 0 1 0 5 】

（ b ）この接続情報 1 1 a を対象にした接続リクエスト以外の各種リクエストが到着すると、図 9 のステップ S 1 0 において接続情報 1 1 a が存在する場合、リクエスト処理の完了時点で追放予定時刻 1 1 b が更新される（ステップ S 1 0 の N o ルート， S 7 ）。すなわち、ステップ S 1 0 の判断の際に、リクエストハンドラを処理するスレッドとクライアント追放用スレッドとの間で接続情報 1 1 a の排他関係が適切にとれていればよい。

20

【 0 1 0 6 】

〔 1 - 6 - 2 〕クライアント側の動作例

次に、クライアント 1 3 0 による動作例を説明する。図 1 1 ~ 図 1 4 は、図 2 に示すクライアント 1 3 0 による接続リクエスト送信処理、リクエスト送信処理、返信受信待ち処理、及び返信受信処理の一例をそれぞれ説明するフローチャートである。

〔 1 - 6 - 2 - 1 〕接続リクエスト送信処理

クライアント 1 3 0 では、接続対象のサーバ 1 0 へ接続を確立するための接続リクエストが送信される。

【 0 1 0 7 】

具体的には、図 1 1 に示すように、送信処理部 3 4 によりサーバ 1 0 へ接続リクエストが送信され（ステップ S 3 1 ）、返信（リプライ）の受信の待ち合わせが行なわれる（ステップ S 3 2 ，図 1 3 のステップ S 5 1 ~ S 6 0 ）。

30

受信処理部 3 2 によりサーバ 1 0 から返信が受信されると、受信処理部 3 2 により、返信の内容が解析され、接続リクエストによりエラーが発生していないか否かが判断される（ステップ S 3 3 ）。

【 0 1 0 8 】

エラーが発生していないと判断された場合（ステップ S 3 3 の Y e s ルート）、管理部 3 3 により、返信から追放予定時刻 3 1 b が取得され（ステップ S 3 4 ）、現在時刻が取得した追放予定時刻 3 1 b を経過しているか否かが判断される（ステップ S 3 5 ）。

現在時刻が追放予定時刻 3 1 b を経過していないと判断された場合（ステップ S 3 5 の N o ルート）、管理部 3 3 により、接続リクエストが成功したと判断され、サーバ 1 0 に対する接続情報 3 1 a が新規に作成される（ステップ S 3 6 ）。また、管理部 3 3 により、作成した接続情報 3 1 a に取得した追放予定時刻 3 1 b が記録され（ステップ S 3 7 ）、処理が正常終了する（ステップ S 3 8 ）。

40

【 0 1 0 9 】

一方、ステップ S 3 3 において、接続リクエストにエラー（接続エラー）が発生していると判断された場合（ステップ S 3 3 の N o ルート）、処理が異常終了する（ステップ S 3 9 ）。なお、ステップ S 3 3 の N o ルート経由で処理が異常終了するのは、接続リクエストの送信元（例えばクライアント 1 3 0 のシステム）に接続リクエストをリトライするか否かの判断を委譲するためである。

【 0 1 1 0 】

50

また、ステップS 3 5において、現在時刻が追放予定時刻3 1 bを経過していると判断された場合にも（ステップS 3 5のY e s ルート）、処理が異常終了する（ステップS 3 9）。なお、現在時刻が受信したばかりの追放予定時刻3 1 bを経過してしまっている場合、サーバ1 0側、クライアント1 3 0側、及びクライアント - サーバ間の通信経路のいずれかに異常がある可能性が高い。そこで、リクエスト送信元に判断を委譲するため、ステップS 3 5のY e s ルート経由でも処理が異常終了するのである。

【0 1 1 1】

〔1 - 6 - 2 - 2〕リクエスト送信処理

接続リクエスト送信処理は、クライアント1 3 0側でサーバ1 0との間の接続を確立するための処理であるため、接続管理の観点からは特殊なリクエストである。以下、接続リクエスト以外の、既に接続状態が存在していることを前提としてクライアント1 3 0から送信される各種リクエスト（以下の説明では、単にリクエストという）に対する処理について説明する。

【0 1 1 2】

具体的には、図1 2に示すように、管理部3 3により、リクエスト送信先のサーバ1 0の情報を含む接続情報3 1 aが検索され、接続情報3 1 aが存在するか否かが判断される（ステップS 4 1）。なお、接続情報3 1 aが検索されるのは、接続が確立されている状態で、ユーザ等の操作によってクライアント1 3 0がアンマウント（u m o u n t）される場合や、エラー等により接続が切断される場合があるためである。

【0 1 1 3】

接続情報3 1 aが存在しないと判断された場合（ステップS 4 1のN o ルート）、当該リクエストの送信は、接続が確立されていないサーバ1 0に対する処理であるため、処理が異常終了する（ステップS 5 0）。

一方、接続情報3 1 aが存在すると判断された場合（ステップS 4 1のY e s ルート）、管理部3 3により、現在時刻と接続情報3 1 aに記録されている追放予定時刻3 1 bとが比較され、現在時刻が追放予定時刻3 1 bを経過しているか否かが判断される（ステップS 4 2）。

【0 1 1 4】

現在時刻が追放予定時刻3 1 bを経過していると判断された場合（ステップS 4 2のY e s ルート）、管理部3 3により、送信処理部3 4へ接続リクエスト送信処理が指示される（ステップS 4 3，図1 1のステップS 3 1～S 3 9）。

次いで、受信処理部3 2により、接続処理が正常終了したか否かが判断され（ステップS 4 4）、異常終了したと判断された場合（ステップS 4 4のN o ルート）、管理部3 3により、接続情報3 1 aが存在するか（作成されたか）否かが判断される（ステップS 4 9）。接続情報3 1 aが存在しないと判断された場合（ステップS 4 9のN o ルート）、接続リトライのために、処理がステップS 4 3に移行する。一方、接続情報3 1 aが存在すると判断された場合（ステップS 4 9のY e s ルート）、処理がステップS 4 2に移行する。

【0 1 1 5】

なお、ステップS 4 9において、接続リクエスト送信が異常終了したにもかかわらず、接続情報3 1 aの存在有無を再度確認する。これは、同時刻に同一クライアント1 3 0から同一サーバ1 0に対して2以上のリクエストが発行され、他のリクエストが先に接続確立を完了していた場合を想定したためである。

一方、ステップS 4 4において、接続処理が正常終了したと判断された場合（ステップS 4 4のY e s ルート）、送信処理部3 4により、取得した追放予定時刻3 1 bに基づきサーバ1 0へリクエストが送信される（ステップS 4 5）。また、受信処理部3 2により、返信受信待ち処理が実行される（ステップS 4 6，図1 3のステップS 5 1～S 6 0）。

【0 1 1 6】

なお、ステップS 4 2において、現在時刻が追放予定時刻3 1 bを経過していないと判

10

20

30

40

50

断された場合（ステップS 4 2のN o ルート）、処理がステップS 4 5に移行する。

次いで、受信処理部3 2により、返信受信待ち処理が正常終了したか否かが判断される（ステップS 4 7）。管理部3 3により異常終了したと判断された場合（ステップS 4 7のN o ルート）、リクエストをリトライすべきか否かを送信元の処理の判断に委ねるため、処理が異常終了する（ステップS 5 0）。

【0 1 1 7】

一方、管理部3 3により返信受信待ち処理が正常終了したと判断された場合（ステップS 4 7のY e s ルート）、処理が正常終了する（ステップS 4 8）。

〔1 - 6 - 2 - 3〕返信受信待ち処理

クライアント1 3 0は、サーバ1 0へ接続リクエスト又は各種リクエストを送信すると、サーバ1 0からのリクエスト返信を待ち合わせる処理を行なう。ただし、リクエストが非同期リクエストであった場合、クライアント1 3 0は、リクエスト返信を待たずに即座に復帰する。

【0 1 1 8】

図1 3に示すように、クライアント1 3 0によりサーバ1 0へ接続リクエスト又は各種リクエストが送信されると、受信処理部3 2により、送信したリクエストが同期通信であるか否かが判断される（ステップS 5 1）。送信したリクエストが同期通信ではない、つまり非同期通信であると判断された場合（ステップS 5 1のN o ルート）、処理が即座に正常終了する（ステップS 5 6）。

【0 1 1 9】

一方、送信したリクエストが同期通信であると判断された場合（ステップS 5 1のY e s ルート）、受信処理部3 2により、リクエスト返信の受信の待ち合わせが行なわれる（ステップS 5 2）。また、受信処理部3 2により、受信されると返信受信処理が行なわれる（ステップS 5 4、図1 4のステップS 6 1～S 7 1）。そして、受信処理部3 2により、返信受信処理が正常終了しているか否かが判断され（ステップS 5 5）、正常終了していると判断された場合（ステップS 5 5のY e s ルート）、リクエスト返信待ち処理も正常終了する（ステップS 5 6）。

【0 1 2 0】

一方、返信受信処理が異常終了していると判断された場合（ステップS 5 5のN o ルート）、受信処理部3 2により、返信受信処理から返されたエラーが、現在時刻が追放予定時刻3 1 bを経過したために発生したエラー以外のエラーであるか否かが判断される（ステップS 5 7）。

エラーが、現在時刻が追放予定時刻3 1 bを経過したために発生したエラー以外のエラーであると判断された場合（ステップS 5 7のY e s ルート）、処理が返信受信処理と同じエラーで異常終了する（ステップS 6 0）。一方、エラーが、現在時刻が追放予定時刻3 1 bを経過したために発生したエラーである場合（ステップS 5 7のN o ルート）、送信処理部3 4により、非同期通信のp i n g リクエストが作成される（ステップS 5 8）。そして、送信処理部3 4により、p i n g リクエストがサーバ1 0へ非同期送信されて（ステップS 5 9、図1 2のステップS 4 1～S 5 0）、処理が正常終了する（ステップS 5 6）。

【0 1 2 1】

なお、ステップS 5 7のN o ルート経由で処理が正常終了する理由は、現在時刻が追放予定時刻3 1 bを超過していることを除いて、全ての処理が正常に完了しているためである。つまり、サーバ1 0に対するリクエストの送信・返信処理自体は正常に完了しており、本来の目的は達成できているためである。

また、ステップS 5 8及びS 5 9で送信処理部3 4がp i n g リクエストをサーバ1 0へ送信するのは、次にクライアント1 3 0がサーバ1 0へリクエストを送信するときに、クライアント追放によるエラーが返されないようにするためである。すなわち、クライアント1 3 0は、ステップS 5 8及びS 5 9において、何らかのリクエスト（この場合、システムが発行するp i n g リクエスト）をサーバ1 0へ送信し、追放復旧処理を実行して

おくのである。

【 0 1 2 2 】

〔 1 - 6 - 2 - 4 〕 返信受信処理

リクエスト返信受信処理は、同期通信の場合、クライアント 1 3 0 におけるリクエストの送信元によって開始される。また、リクエスト返信受信処理は、非同期通信の場合、例えばリクエスト送受信を管理するリクエスト送信元とは別のスレッドからコールバック関数の呼び出し等によって開始される。例えば、クライアント 1 3 0 は、サーバ 1 0 からリクエストに対する返信を受信すると、リクエスト固有の返信受信処理を行なう。

【 0 1 2 3 】

具体的には、図 1 4 に示すように、受信処理部 3 2 により、受信したリクエスト返信メッセージの内容が解析され（ステップ S 6 1 ）、通信エラーが発生しているか否かが判断される（ステップ S 6 2 ）。

受信処理部 3 2 により、通信エラーが発生したと判断された場合（ステップ S 6 2 の Y e s ルート）、追放予定時刻 3 1 b が更新されておらず、処理が即座に異常終了する（ステップ S 7 1 ）。なお、ステップ S 6 2 の Y e s ルートに進む場合としては、サーバ 1 0 へリクエストが届いていない場合も含まれる。

【 0 1 2 4 】

一方、通信エラーが発生していないと判断された場合（ステップ S 6 2 の N o ルート）、受信処理部 3 2 により、各リクエストに固有の返信受信処理が実行される（ステップ S 6 3 ）。そして、管理部 3 3 により、返信メッセージ中の追放予定時刻 3 1 b が取得され（ステップ S 6 4 ）、取得された追放予定時刻 3 1 b が接続情報 3 1 a に上書き保存される（ステップ S 6 5 ）。

【 0 1 2 5 】

また、受信処理部 3 2 により、リクエスト固有の返信受信処理がエラー終了したか否かが判断される（ステップ S 6 6 ）。正常終了したと判断された場合（ステップ S 6 6 の N o ルート）、管理部 3 3 により、現在時刻が新しい追放予定時刻 3 1 b を経過したか否かが判断される（ステップ S 6 7 ）。現在時刻が新しい追放予定時刻 3 1 b を経過していないと判断された場合（ステップ S 6 7 の N o ルート）、処理が正常終了する（ステップ S 6 8 ）。なお、現在時刻が新しい追放予定時刻 3 1 b を経過したと判断された場合（ステップ S 6 7 の Y e s ルート）、サーバ 1 0 側でクライアント 1 3 0 が追放されたため、処理が異常終了する（ステップ S 7 0 ）。

【 0 1 2 6 】

一方、ステップ S 6 6 において、異常終了したと判断された場合（ステップ S 6 6 の Y e s ルート）、管理部 3 3 により、現在時刻が新しい追放予定時刻 3 1 b を経過したか否かが判断される（ステップ S 6 9 ）。現在時刻が新しい追放予定時刻 3 1 b を経過していないと判断された場合（ステップ S 6 9 の N o ルート）、処理が異常終了する（ステップ S 7 1 ）。なお、現在時刻が新しい追放予定時刻 3 1 b を経過したと判断された場合（ステップ S 6 9 の Y e s ルート）、サーバ 1 0 側でクライアント 1 3 0 が追放されたため、処理が異常終了する（ステップ S 7 0 ）。

【 0 1 2 7 】

〔 2 〕 その他

以上、本発明の好ましい実施形態について詳述したが、本発明は、係る特定の実施形態に限定されるものではなく、本発明の趣旨を逸脱しない範囲内において、種々の変形、変更して実施することができる。

例えば、上述した説明では、1つのクライアント 1 3 0 がサーバ 1 0 と通信を行なう場合の分散ファイルシステム 1 の動作について説明したが、図 2 に示すように、複数のクライアント 1 3 0 がサーバ 1 0 と通信を行なう場合も同様である。この場合、図 1 5 に示す例のように、クライアント A は、処理 T 1 1 0 でサーバから分散ロックを付与されると、定期的に（例えば 2 5 秒よりも短い間隔で）サーバへリクエストを送信することで、追放予定時刻 3 1 b を更新することができる。なお、クライアント B は、処理 T 1 2 0 でサー

10

20

30

40

50

バへ分散ロック要求を発行するが、サーバから返信が返っていないため、処理 T 1 6 0 まではサーバにおいてクライアント B の追放予定時刻 3 1 b の計算は行なわれない。

【 0 1 2 8 】

また、サーバ 1 0 及びクライアント 1 3 0 は、同一の追放予定時刻 1 1 b 及び 3 1 b を保持するものとして説明したが、これに限定されるものではない。例えば、サーバ 1 0 は、クライアント 1 3 0 - サーバ 1 0 間のネットワークの遅延やクライアント 1 3 0 及びサーバ 1 0 での処理遅延等を考慮して、サーバ 1 0 が保持する追放予定時刻 1 1 b よりも早い時刻を示す追放予定時刻 3 1 b をクライアント 1 3 0 へ通知してもよい。

【 0 1 2 9 】

なお、一実施形態に係るサーバ 1 0 及びクライアント 1 3 0 の各種機能の全部もしくは一部は、コンピュータ (CPU, 情報処理装置, 各種端末を含む) が所定のプログラムを実行することによって実現されてもよい。

そのプログラムは、例えばフレキシブルディスク、CD、DVD、ブルーレイディスク等のコンピュータ読取可能な記録媒体 (例えば図 3 に示す記録媒体 1 0 h) に記録された形態で提供される。なお、CD としては、CD-ROM、CD-R、CD-RW 等が挙げられる。また、DVD としては、DVD-ROM、DVD-RAM、DVD-R、DVD-RW、DVD+R、DVD+RW 等が挙げられる。この場合、コンピュータはその記録媒体からプログラムを読み取って内部記憶装置または外部記憶装置に転送し格納して用いる。

【 0 1 3 0 】

〔 3 〕 付記

以上の実施形態に関し、更に以下の付記を開示する。

(付記 1)

情報処理装置と、前記情報処理装置との間で確立された接続を用いて前記情報処理装置と通信を行なう端末装置とを有する情報処理システムにおいて、

前記情報処理装置は、

前記接続を解除する予定時刻を前記端末装置に通知し、

前記端末装置は、

前記情報処理装置へ要求を送信する際に、現在時刻が前記情報処理装置から通知された予定時刻を経過しているか否かを判断する判断部と、

前記判断部により前記現在時刻が前記予定時刻を経過していると判断された場合、前記情報処理装置へ前記要求を送信する前に、前記情報処理装置との間で接続を確立するための接続要求を前記情報処理装置へ送信する送信部と、を有することを特徴とする、情報処理システム。

【 0 1 3 1 】

(付記 2)

前記情報処理装置は、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断し、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除する接続解除部、を有する

ことを特徴とする、付記 1 記載の情報処理システム。

【 0 1 3 2 】

(付記 3)

前記情報処理装置は、

前記予定時刻を前記端末装置に関する接続情報に対応付けて管理する接続管理部、をさらに有し、

前記接続解除部は、

前記現在時刻が前記接続管理部の管理する前記予定時刻を経過していると判断した場合、前記接続情報を無効化する

ことを特徴とする、付記 2 記載の情報処理システム。

10

20

30

40

50

【 0 1 3 3 】

(付 記 4)

前記接続管理部は、

前記端末装置から接続要求又は前記接続要求以外の要求を受けると、現在時刻に所定時間を加算して予定時刻を取得し、取得した前記予定時刻を前記端末装置に関する接続情報に対応付けて管理する

ことを特徴とする、付記 3 記載の情報処理システム。

【 0 1 3 4 】

(付 記 5)

前記情報処理装置は、

前記端末装置から接続要求又は前記接続要求以外の要求を受けると、前記接続要求又は前記要求への応答に前記予定時刻を含めて、前記端末装置へ通知する通知部、を有することを特徴とする、付記 1 ～ 4 のいずれか 1 項記載の情報処理システム。

【 0 1 3 5 】

(付 記 6)

前記端末装置は、

前記情報処理装置から受信した予定時刻を前記情報処理装置に関する接続情報に対応付けて管理する時刻管理部と、

前記判断部により前記現在時刻が前記予定時刻を経過していると判断された場合、前記接続情報を無効化する接続情報管理部と、をさらに有する

ことを特徴とする、付記 1 ～ 5 のいずれか 1 項記載の情報処理システム。

【 0 1 3 6 】

(付 記 7)

前記判断部は、

前記情報処理装置へ要求を送信する際に、現在時刻が前記時刻管理部により管理される前記予定時刻を経過しているか否かを判断し、

前記送信部は、

前記判断部により前記現在時刻が前記予定時刻を経過していないと判断された場合、前記接続情報に基づいて前記情報処理装置へ前記要求を送信する一方、前記判断部により前記現在時刻が前記予定時刻を経過していると判断された場合、前記接続情報管理部により無効化された接続情報を用いずに、前記情報処理装置へ前記要求を送信する前に前記情報処理装置へ接続要求を送信する

ことを特徴とする、付記 6 記載の情報処理システム。

【 0 1 3 7 】

(付 記 8)

端末装置との間で確立した接続を用いて前記端末装置と通信を行なう情報処理装置において、

前記端末装置から接続要求を受けると、前記端末装置との間で接続を確立する接続処理を行なう接続管理部と、

前記接続処理により確立する接続を切り離す予定時刻を前記端末装置に通知する通知部と、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断し、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除する接続解除部、を有する

ことを特徴とする、情報処理装置。

【 0 1 3 8 】

(付 記 9)

前記接続管理部は、

前記予定時刻を前記端末装置に関する接続情報に対応付けて管理し、

前記接続解除部は、

前記現在時刻が前記接続管理部の管理する前記予定時刻を経過していると判断した場合、前記接続情報を無効化する
ことを特徴とする、付記 8 記載の情報処理装置。

【 0 1 3 9 】

(付記 1 0)

前記接続管理部は、

前記端末装置から接続要求又は前記接続要求以外の要求を受けると、現在時刻に所定時間を加算して予定時刻を取得し、取得した前記予定時刻を前記端末装置に関する接続情報に対応付けて管理する

ことを特徴とする、付記 8 又は付記 9 記載の情報処理装置。

10

【 0 1 4 0 】

(付記 1 1)

前記通知部は、

前記端末装置から接続要求又は前記接続要求以外の要求を受けると、前記接続要求又は前記要求への応答に前記予定時刻を含めて、前記端末装置へ通知する

ことを特徴とする、付記 8 ~ 1 0 のいずれか 1 項記載の情報処理装置。

【 0 1 4 1 】

(付記 1 2)

端末装置との間で確立した接続を用いて前記端末装置と通信を行なう情報処理装置の制御プログラムにおいて、

20

前記情報処理装置に、

前記端末装置から接続要求を受けると、前記端末装置との間で接続を確立する接続処理を行なわせ、

前記接続処理により確立する接続を切り離す予定時刻を前記端末装置に通知させ、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断させ、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除させる

ことを特徴とする、情報処理装置の制御プログラム。

【 0 1 4 2 】

(付記 1 3)

前記情報処理装置に、

前記予定時刻を前記端末装置に関する接続情報に対応付けて管理させ、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断させ、前記現在時刻が前記予定時刻を経過していると判断した場合、前記接続情報を無効化させる

ことを特徴とする、付記 1 2 記載の情報処理装置の制御プログラム。

30

【 0 1 4 3 】

(付記 1 4)

前記情報処理装置に、

前記端末装置から接続要求又は前記接続要求以外の要求を受けると、前記接続要求又は前記要求への応答に前記予定時刻を含めて、前記端末装置へ通知させる

ことを特徴とする、付記 1 2 又は付記 1 3 記載の情報処理装置の制御プログラム。

40

【 0 1 4 4 】

(付記 1 5)

情報処理装置と、前記情報処理装置との間で確立された接続を用いて前記情報処理装置と通信を行なう端末装置とを有する情報処理システムの制御方法において、

前記情報処理装置が、

前記接続を解除する予定時刻を前記端末装置に通知し、

前記端末装置が、

前記情報処理装置へ要求を送信する際に、現在時刻が前記情報処理装置から通知された予定時刻を経過しているか否かを判断し、

50

前記現在時刻が前記予定時刻を経過していると判断した場合、前記情報処理装置へ前記要求を送信する前に、前記情報処理装置との間で接続を確立するための接続要求を前記情報処理装置へ送信する

ことを特徴とする、情報処理システムの制御方法。

【0145】

(付記16)

前記情報処理装置が、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断し、前記現在時刻が前記予定時刻を経過していると判断した場合、前記端末装置との間で確立された接続を解除する

10

ことを特徴とする、付記15記載の情報処理システムの制御方法。

【0146】

(付記17)

前記情報処理装置が、

前記予定時刻を前記端末装置に関する接続情報に対応付けて管理し、

所定時間ごとに、現在時刻が前記予定時刻を経過しているか否かを判断し、前記現在時刻が前記接続管理部の管理する前記予定時刻を経過していると判断した場合、前記接続情報を無効化する

ことを特徴とする、付記16記載の情報処理システムの制御方法。

【0147】

20

(付記18)

前記情報処理装置が、

前記端末装置から接続要求又は前記接続要求以外の要求を受けると、前記接続要求又は前記要求への応答に前記予定時刻を含めて、前記端末装置へ通知する

ことを特徴とする、付記15～17のいずれか1項記載の情報処理システムの制御方法。

【0148】

(付記19)

前記端末装置が、

前記情報処理装置から受信した予定時刻を前記情報処理装置に関する接続情報に対応付けて管理し、

30

前記現在時刻が前記予定時刻を経過していると判断した場合、前記接続情報を無効化する

ことを特徴とする、付記15～18のいずれか1項記載の情報処理システムの制御方法。

【0149】

(付記20)

前記判断部が、

前記情報処理装置へ要求を送信する際に、現在時刻が前記予定時刻を経過しているか否かを判断し、

前記現在時刻が前記予定時刻を経過していないと判断した場合、前記接続情報に基づいて前記情報処理装置へ前記要求を送信し、

40

前記現在時刻が前記予定時刻を経過していると判断した場合、無効化した接続情報を用いずに、前記情報処理装置へ前記要求を送信する前に前記情報処理装置へ接続要求を送信する

ことを特徴とする、付記19記載の情報処理システムの制御方法。

【符号の説明】

【0150】

1 分散ファイルシステム(情報処理システム)

10, 10-1 MGS(サーバ, 情報処理装置)

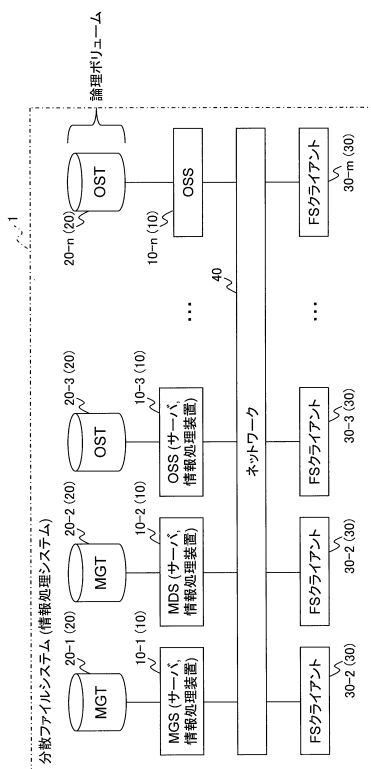
10, 10-2 MDS(サーバ, 情報処理装置)

10, 10-3～10-n OSS(サーバ, 情報処理装置)

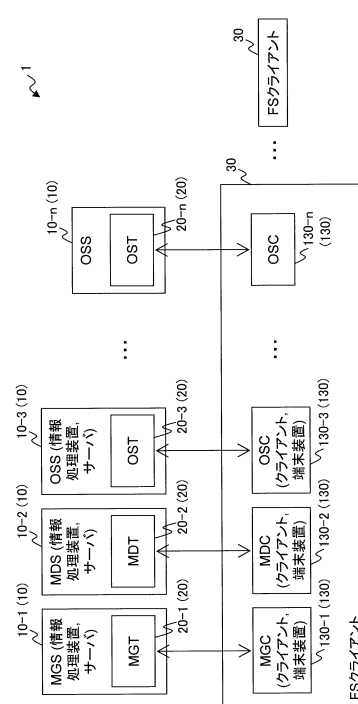
50

10a, 30a	CPU	
10b, 30b	メモリ	
10c, 30c	記憶部	
10d, 30d	ネットワークインタフェース	
10e, 30e	入出力部	
10f, 10h, 30f, 30h	記録媒体	
10g, 30g	読取部	
11, 31, 31-1~31-3	保持部	
11a, 11a-1~11a-3, 31a, 31a-1~31a-3	接続情報	
11b, 11b-1~11b-3, 31b, 31b-1~31b-3	追放予定時刻	10
(予定時刻)		
12, 32	受信処理部	
13	リクエスト処理部(接続管理部, 通知部)	
14	追放処理部(接続解除部)	
20, 20-1	MGT(論理ボリューム)	
20, 20-2	MDT(論理ボリューム)	
20, 20-3~20-n	OST(論理ボリューム)	
30, 30-1~30-m	FSクライアント	
33	管理部(時刻管理部, 判断部, 接続情報管理部)	
34	送信処理部(送信部)	20
40	ネットワーク	
130, 130-1	MGC(クライアント, 端末装置)	
130, 130-2	MDC(クライアント, 端末装置)	
130, 130-3~130-n	OSC(クライアント, 端末装置)	

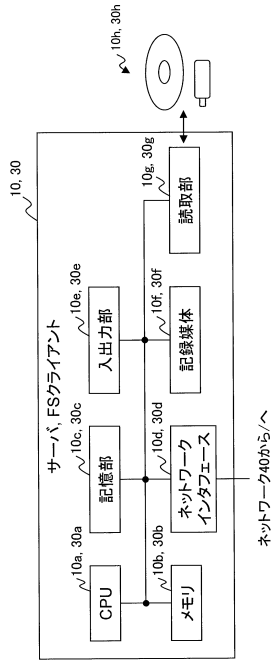
【図1】



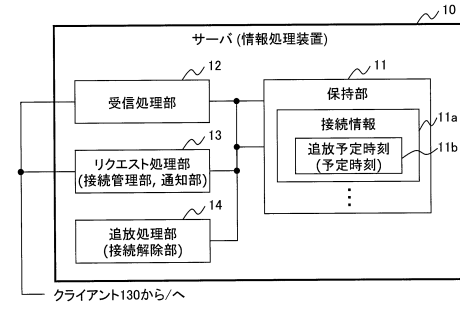
【図2】



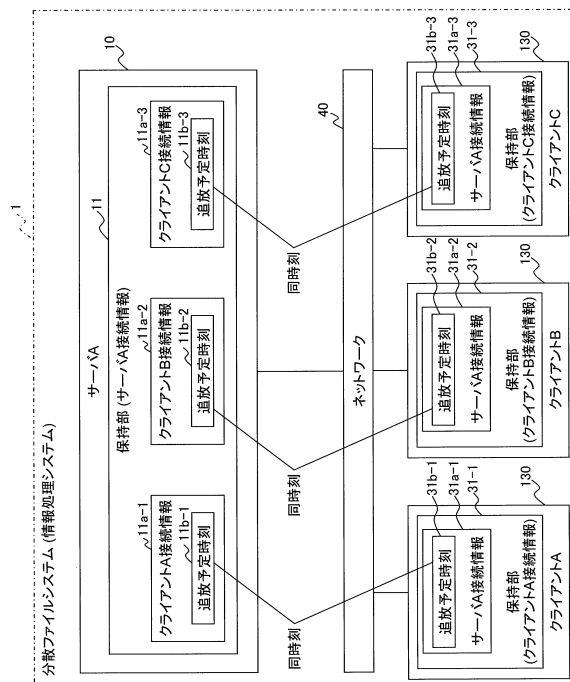
【図 3】



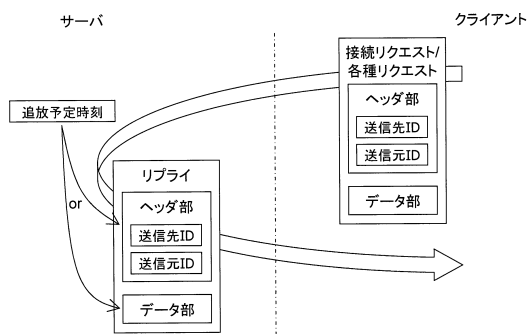
【図 4】



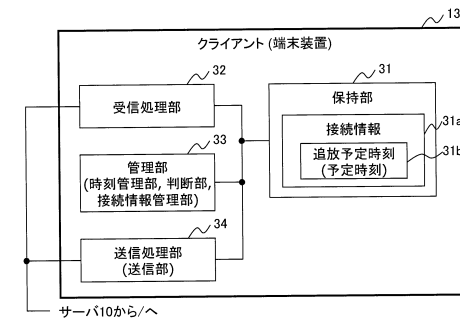
【図 5】



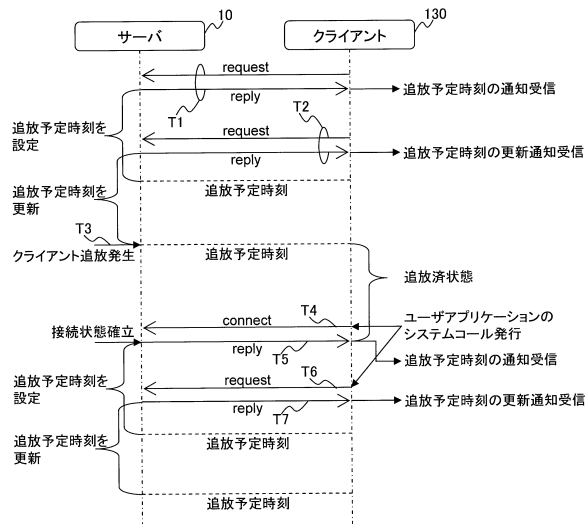
【図 6】



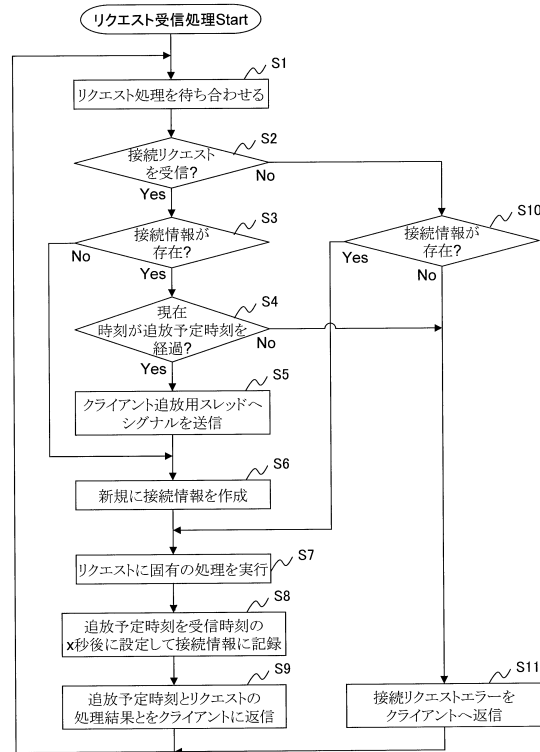
【図 7】



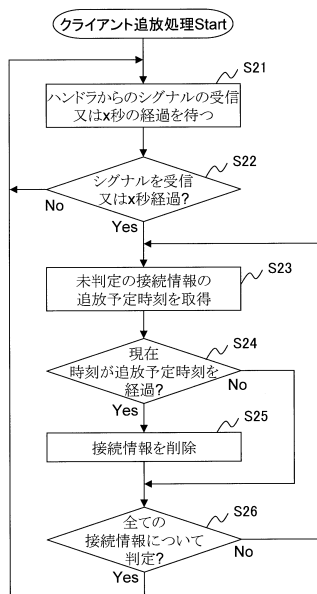
【図 8】



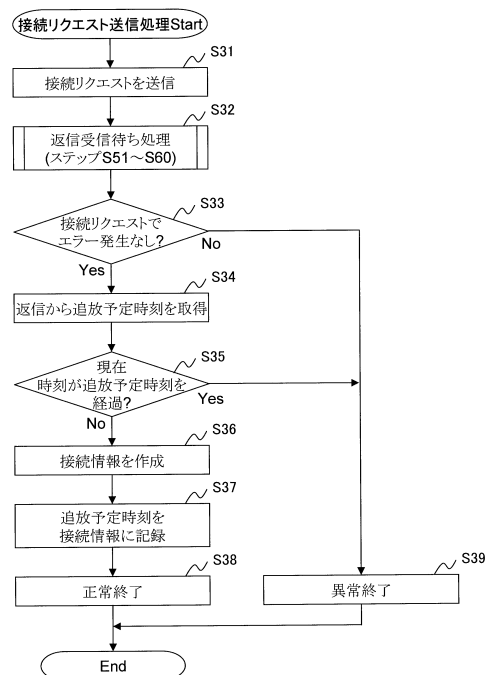
【図 9】



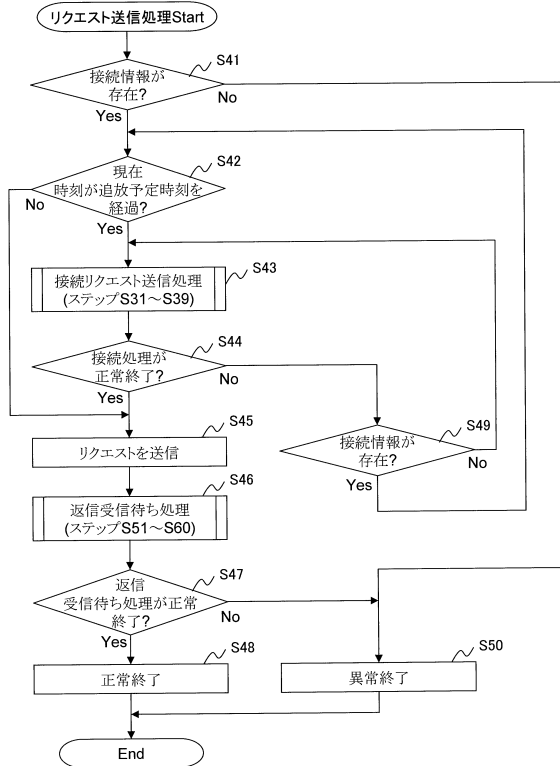
【図 10】



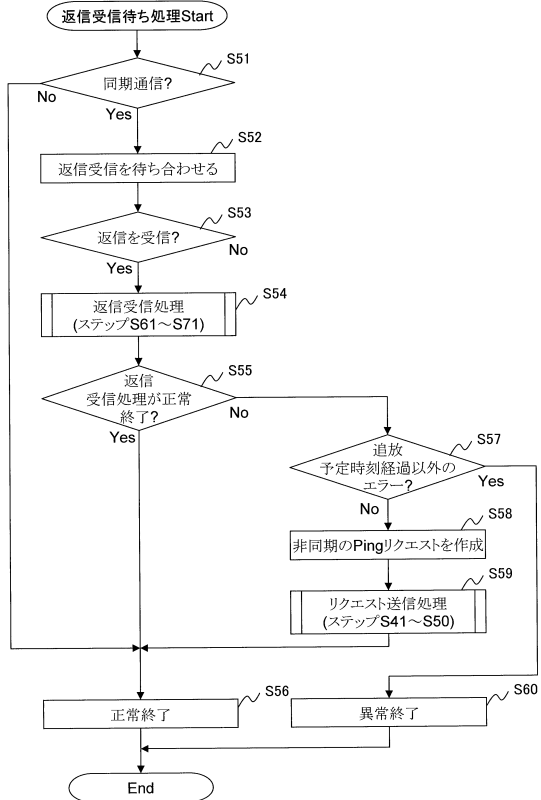
【図 11】



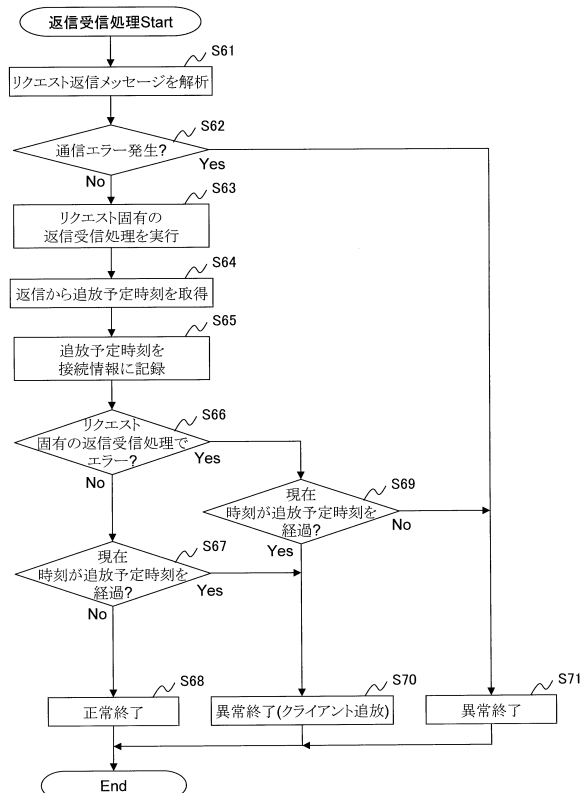
【図 12】



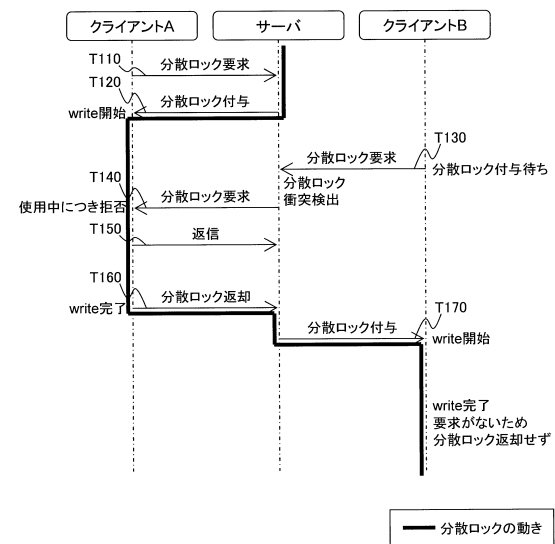
【図 13】



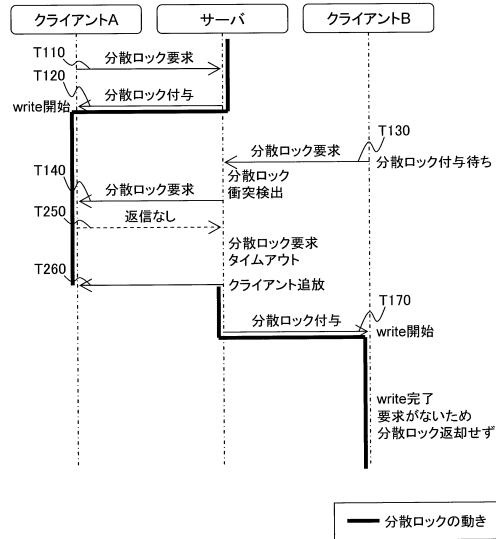
【図 14】



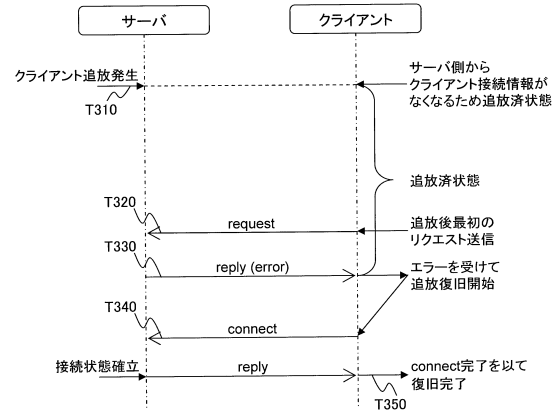
【図 15】



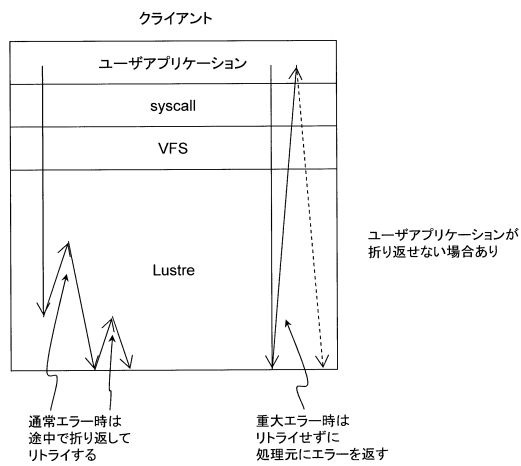
【図 16】



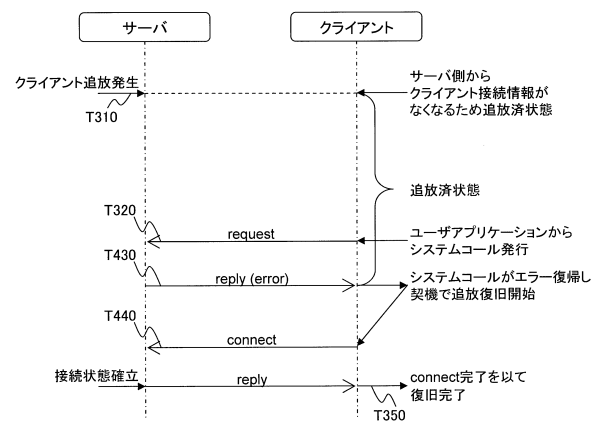
【図 17】



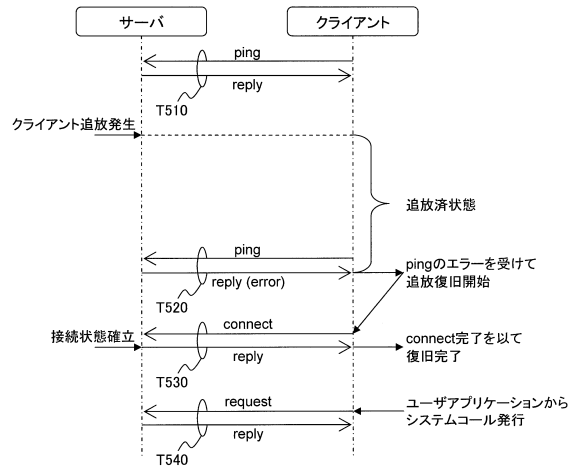
【図 18】



【図 19】



【図 20】



フロントページの続き

(56)参考文献 特開2005-346573(JP,A)
特開2009-48510(JP,A)
特開2000-224260(JP,A)
特開平10-164054(JP,A)
特開2007-36624(JP,A)
米国特許出願公開第2011/0119241(US,A1)
米国特許出願公開第2008/0140941(US,A1)
米国特許出願公開第2005/0044240(US,A1)
米国特許第7406523(US,B1)

(58)調査した分野(Int.Cl., DB名)
G06F 15/00
G06F 12/00
G06F 13/00