

分散ファイルシステムにおける冗長化方式の一検討

A Study on Server Redundancy for Distributed File System

西谷 明彦†
Akihiko NISHITANI寺下 雅人†
Masahito TERASHITA大岸 智彦†
Tomohiko OGISHI

1. まえがき

クラウドサービスの普及に伴い、多くのシステムの基盤として、PC が持つストレージを統合的に利用する分散ファイルシステムの活用が着目されている。分散ファイルシステムでは、低価格な汎用 PC による構築が可能ことや、ユーザ数等の増加に応じて必要な分だけハードディスクや CPU などのリソースを増やせるスケールアウト性を有するため、設備導入コストの削減への寄与が期待される一方、数台の PC に障害が発生した場合に、保管するデータを安全に保護し、システム全体の障害を防ぐことで、サービスが停止せずに運用可能とするための可用性・信頼性も求められている。

ハードディスクなどの突発的な故障に起因するシステム障害を防ぐには、単一障害点 (SPoF: Single Point of Failure) の除去が重要である。筆者らは、既存の分散ファイルシステムの仕組みを生かしつつ、単一障害点となりうるメタデータサーバを冗長化する手法を考案した。本稿ではその概要について述べる。

2. 現状の構成と課題

現状の分散ファイルシステムとして、Gfarm[1]や Lustre[2]など、ユーザが作成したファイルを保管する多数のファイルサーバ、各ファイルサーバが格納するファイルのメタデータを一元管理するメタデータサーバ、ユーザアプリケーションに対して分散ファイルシステムへのアクセス手段を提供するプロキシサーバで構成されたアーキテクチャを仮定する。この構成では、1 台のメタデータサーバが、全てのユーザからの要求を処理しデータの一貫性を維持するため、メタデータサーバ上でプロセスの異常動作などのソフトウェア障害や、ハードディスク故障などのハードウェア障害が発生した場合に、システム全体の停止を招くこととなる。これを回避するため、Act-HotStandby 構成でメタデータサーバを冗長化する手法が考えられるが、ユーザ要求処理と独立にメタデータの同期を行うため、サービスのダウンタイム発生を防ぐことが難しいという問題がある。

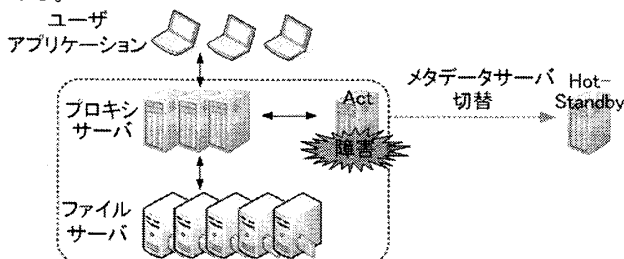


図1 現状の分散ファイルシステムの課題

†株式会社 KDDI 研究所, KDDI R&D Laboratories, Inc.

3. 提案する冗長化方式

既存の分散ファイルシステムは、メタデータの管理や、遅延や CPU 使用率などをもとにファイルサーバを選択する仕組みなど、特徴的な機能を有する。既存の分散ファイルシステム上に冗長化機能を導入することで、既存の分散ファイルシステムの特徴を生かしつつ、課題となっているメタデータサーバの冗長化を行う方式を考案した。

3.1 アーキテクチャ

提案する冗長化方式のアーキテクチャを図2に示す。本構成では、既存の分散ファイルシステム相当のシステムを複数システム導入する。ユーザからの要求に対して、全てのシステムが同時に動作することで、いずれかのシステムで障害が発生しても他のシステムにてサービスを継続する。DNS サーバは、ユーザからの問い合わせに対して、最初に要求を受け付けるシステムを選択する仕組みを提供する。生存監視サーバは、各システム内のサーバの故障や復帰を迅速に検知する機能を有し、DNS 登録情報と連動させてシステム別の死活を制御する役割を果たす。

なお、個々のシステム内では、プロキシサーバ、メタデータサーバはそれぞれ1台ずつ存在し、ファイルサーバにおいては、ファイルの複製は行わない構成を想定している。

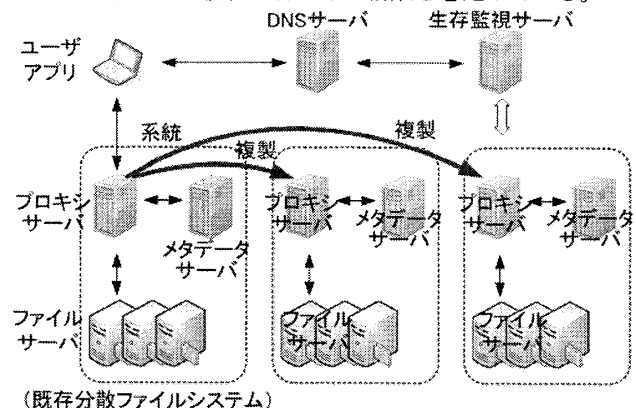


図2 提案するアーキテクチャ

3.2 正常手順

提案方式における正常手順を以下に示す。ファイル書き込みの場合、ユーザからのファイルアクセス要求を DNS サーバが選択した一次プロキシサーバが受け取り、他の二次プロキシサーバに要求を複製することで各システムにおいて並列に同一の要求処理を実施する。二次プロキシサーバより一定数の正常応答が得られたとき、ユーザに結果を通知する(図3参照)。本正常応答数については、応答性能を重視する際は1に近い値とし、データの可用性を重視する際はシステム数-1に近い値とする。読み込みの場合、一次プロキシサーバ自身での読み込みが失敗した場合に限り、他の二

次プロキシサーバに読み込み要求を複製し、その結果を持ってユーザ端末に応答する。

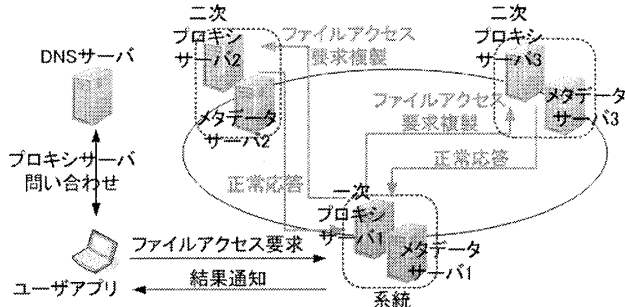


図3 ファイルアクセス命令の複製

3.3 データの一貫性保証

ファイル書き込み要求などにおいて、いずれかのシステムで失敗あるいは著しい応答遅延が発生した場合には、各システムのメタデータサーバが保持するメタデータに差分が生じる。このとき、図4に従いデータの一貫性保証を行う仕組みを実現する。

一次プロキシサーバは、二次プロキシサーバからエラー応答を受信した際、監視サーバにエラー報告を行う。監視サーバは、各プロキシサーバからのエラー履歴をもとにシステム間のファイルの差分を補完するため、各プロキシサーバに対し、不足しているファイルを相互に転送するように指示する。このように、システム間のメタデータおよびファイルの同期は、ユーザの要求処理とは独立なタイミングで行う。

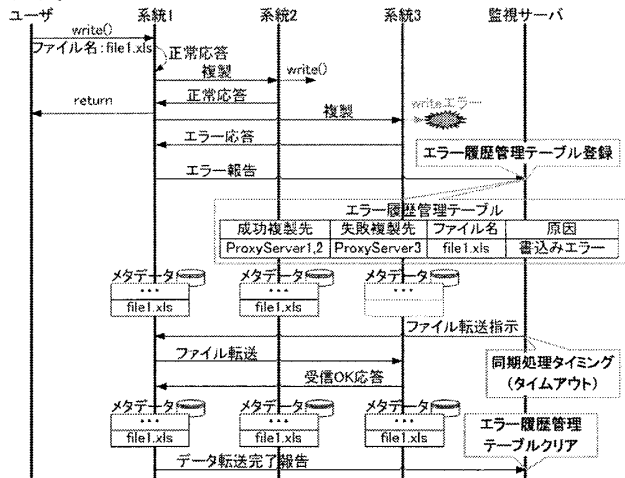


図4 データの一貫性保証

3.4 システム別の死活監視

提案の構成において、各システムのプロキシサーバに障害がある場合、ユーザ要求の失敗という重度の障害が発生する。これを避けるため、監視サーバではシステムごとの死活監視を行い、プロキシサーバが故障したシステムをDNSサーバが選択しないような仕組みを設ける。この仕組みについて図5を用いて説明する。

全てのシステムが正常な場合、DNSサーバは、システム1~3のいずれのプロキシサーバも選択する可能性がある。監視サーバは、ping およびプロセス稼働確認により、定期的、各システムのプロキシサーバの死活監視を行う。応答が無い、あるいはプロセスが稼働していない場合、プロキシサーバ

の故障と判定し、DNSサーバ上にて、そのシステムのプロキシサーバの登録解除手続きを行う。死活監視による故障が確認できないことも考慮し、3.3節で述べたエラー履歴も故障検知の材料として利用する。特定のシステムのエラーが多い場合は、そのシステムの故障と判断し、同様にDNS登録解除を行う。監視サーバは、故障したシステムが正常復帰した時点でDNSサーバにそのシステムのプロキシサーバを再登録する。

なお、このシステムの切り離し・復旧処理については故障時の対応を想定したものであり、システム構成変更による恒久的なシステム数の変更については対象としていない。システム復旧後、データの一元性保証のための同期処理を行う。

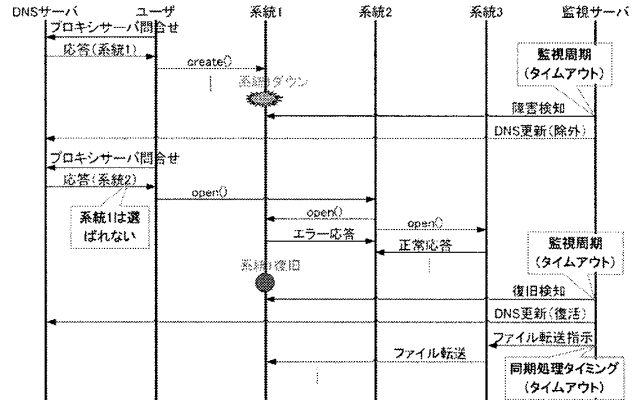


図5 システム別の死活監視

3.5 考察

図2では、プロキシサーバ、メタデータサーバが、それぞれ別々のマシンで動作する例を示したが、サーバプログラムを1台のマシン上で動作させることによりリソースの効率的利用が可能である。また、ファイルサーバについては、異なるシステムのファイルサーバを同一マシン上で動作することも可能であり、これによりディスク領域の効率的利用にも貢献できると考えられる。

一方で、今回の設計では、ユーザ数が増えた場合に、プロキシサーバ/メタデータサーバの処理能力を維持するためにシステム数を増加させる必要がある。全てのシステムが並列に動作することを前提としているため、スケーラビリティに課題が残る。

4. まとめ

本稿では、既存の分散ファイルシステムにおける、メタデータサーバが単一障害点となり得る問題に対し、その対策としての冗長化手法と設計について述べた。今後の課題として、スケーラビリティの向上が挙げられる。

謝辞

日頃よりご指導頂く、株式会社 KDDI 研究所秋葉所長、中島副所長、及び鈴木取締役様に感謝致します。

参考文献

- [1] 建部 修見, 曾田 哲之, "広域分散ファイルシステム Gfarm v2 の実装と評価" 情報処理学会研究報告, 2007-HPC-113, 2007.
- [2] http://wiki.lustre.org/index.php/Main_Page