

# Gfarmにおけるクライアントキャッシュの効率化手法の提案

石 黒 駿<sup>†1</sup> 大 山 恵 弘<sup>†1,†2</sup>

## 1. はじめに

大規模なデータを扱うために、分散ファイルシステムが注目されている。分散ファイルシステムとは、複数のノードにデータを分散して配置し、それらのノードをまとめて一つのストレージとして提供するファイルシステムである。分散ファイルシステムの一つとして、Gfarm<sup>1</sup>が挙げられる。Gfarmは広域ネットワークで利用可能な分散ファイルシステムであり、高容量、高信頼、高性能を実現している。Gfarmはメタデータサーバ及びI/Oサーバから構成される。メタデータサーバは、Gfarmファイルシステム上のファイルのメタデータを一括して管理する。I/Oサーバは複数存在し、ファイルの内容を保存する。GfarmのクライアントがGfarmファイルシステム上のファイルにアクセスする際は、まずメタデータサーバへアクセスし、どのI/Oサーバに目的のファイルがあるかを知り、次にI/Oサーバにアクセスしてデータをやりとりする。

クライアントは、ユーザレベルファイルシステムを実現するためのフレームワークであるFUSE<sup>2</sup>を利用して、Gfarmファイルシステムをマウントできる。通常Gfarmファイルシステムにアクセスするためには、クライアントは専用のAPIを用いる必要があるが、マウントすることにより、システムコールを利用してGfarmファイルシステムにアクセスできる。FUSEにはカーネルモジュールが用いられているため、ユーザレベルファイルシステムはカーネルのキャッシュを利用できる。したがって、マウントされたGfarmファイルシステム上のファイル内容も、クライアントカーネルのページキャッシュにキャッシュされる。しかしながら、現在の実装では、Gfarmファイルシステム上のファイル内容とページキャッシュの内容の一貫性を保つために、ファイルをopenする度にページキャッシュは破棄される。Gfarmファイルシステム上のファイルの更新は、クライアントには通知されず、ファイルオープン時にキャッシュが残っていると、古いデータにアクセスすることになってしまうためである。ゆ

えに、Gfarmファイルシステム上のファイル内容とページキャッシュの内容が一致し本来ページキャッシュを破棄する必要がない場合にもキャッシュが破棄されるため、性能を損なっている。

そこで本研究では、ファイルopen時にページキャッシュの内容が最新であり破棄する必要がない場合に、ページキャッシュを保持するための手法を提案する。これによりページキャッシュの内容が最新でないときのみキャッシュを破棄するため、キャッシュを有効活用できる。

## 2. 方 針

FUSEは図1に示すように、カーネルモジュールであるFUSEモジュールとユーザレベルで動作するFUSEデーモンを連携させることにより、ユーザレベルファイルシステムを実現する。FUSEモジュールが、ユーザレベルファイルシステムに対するリクエストを受け取り、実際の処理をFUSEデーモンに要求する。FUSEデーモンは実際の処理を完了すると結果をFUSEモジュールへ返し、リクエストが完了する。FUSEデーモンは多くの場合、ネットワークで通信したりext3などのローカルファイルシステムにアクセスしたりする。

Gfarmファイルシステムのマウントでは、FUSEデーモンとしてgfarm2fsが利用される。gfarm2fsは、FUSEモジュールからリクエストを受け取ると、メタデータサーバやI/Oサーバとやりとりして、結果をFUSEモジュールへ返す。提案手法の方針として、gfarm2fsがopenリクエスト処理時にそのファイルのページキャッシュの内容が最新かどうかをチェックし、最新の場合はキャッシュを破棄しないようにFUSEモジュールへ伝える。

## 3. 提案手法

まず、ページキャッシュの内容が最新であるかどうかを調べる方法であるが、これにはGfarmのinode番号と世代番号を利用する。inode番号は、ファイル固有の一意な値であり、世代番号はそのinodeが表すファイルが更新されるとインクリメントされる値である。提案手法では、まずopenしたファイルのinode番号とそのinodeの世代番号をペアで覚えておく。次の

†1 電気通信大学

The University of Electro-Communications

†2 独立行政法人科学技術振興機構, CREST

JST, CREST

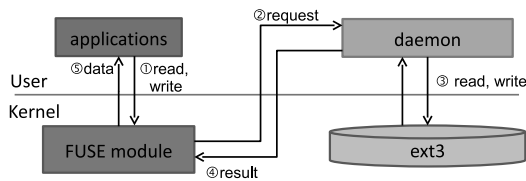


図 1 FUSE のアーキテクチャ

open 時に open 対象のファイルの inode 番号とその世代番号をメタデータサーバから取得し、その inode 番号を元に覚えておいた世代番号と取得した世代番号を比較する。もし、世代番号が一致していればそのファイルは更新されていないため、ページキャッシュを破棄せず、そうでない場合は破棄する。

次に gfarm2fs から FUSE モジュールにページキャッシュを破棄しないよう伝える方法である。FUSE では FUSE デーモンを起動する際のオプションで、ファイル open 時もキャッシュを破棄しないように指定することができるが、分散ファイルシステムのようなファイルシステムでは、ファイル内容の一貫性の問題からこのオプションは利用されない。FUSE デーモンを開発するためのライブラリは、ファイルの open 処理時に、このオプションの値を調べ、ページキャッシュを破棄するか否かを FUSE モジュールへ伝える。FUSE ライブラリは、ファイル情報を持つ構造体のメンバに値を設定することで、ページキャッシュを破棄するかどうかの情報を FUSE モジュールへ伝えている。そして、このメンバには FUSE デーモンからアクセスすることができる。そこで提案手法では、このメンバの値を FUSE ライブラリが設定した後に上書きすることにより、ページキャッシュの破棄操作を制御する。提案手法は、FUSE デーモンの変更のみで実現できる。

#### 4. 評価

提案手法の有効性を調べるために、評価実験を行った。実験では、1GB のファイルを Gfarm ファイルシステム上に用意し、このファイルを 4KB ずつ read するプログラムの実行時間を計測した。キャッシュの効果を確かめるために、プログラムを 2 回連続で実行し、2 回目以降の実行時間を計測した。実験は、提案手法を用いない場合と用いる場合に対して行った。実験環境は、メタデータサーバ 1 台、I/O サーバ 1 台、クライアントのノード 1 台の計 3 ノードであり、クライアントのノードに Gfarm ファイルシステムをマウントした。すべてのノードの性能は、CPU が Intel Xeon 2.40GHz × 2、メモリが 48GB、HDD が 15,000rpm 600GB であり、ノード間は InfiniBand にて接続されている。また OS は Cent OS 5.5 64bit (kernel-2.6.18)、Gfarm のバージョンは 2.4.2、FUSE のバージョンは 2.7.4、

表 1 Gfarm ファイルシステム上の 1GB のファイルを 4KB ずつ read するプログラムの実行時間

|          | 提案手法なし | 提案手法あり |
|----------|--------|--------|
| 実行時間 (s) | 2.96   | 0.330  |

Gfarm2fs のバージョンは 1.2.3 である。実験の結果を表 1 に示す。

実験の結果、ページキャッシュの内容が最新のファイルに対する read では、クライアントのページキャッシュが利用されるため、性能が向上することが確認できた。提案手法により、場合に応じてページキャッシュを破棄するかしないかを open 時に指定ことができ、より効率的にページキャッシュを利用できる。

#### 5. 関連研究

Gfarm を対象にした FUSE のキャッシュ機構の研究として、石黒らの研究<sup>3)</sup>がある。この研究では、同一の端末で複数のユーザが Gfarm ファイルシステムをそれぞれ別のディレクトリにマウントした場合に、本来同一ファイルのキャッシュが複数できてしまうものを、一つにまとめている。

#### 6. 現状と今後

マウントした Gfarm ファイルシステムにおいて、ページキャッシュを open 時に破棄するかしないかを制御する手法を提案した。これにより、ページキャッシュの内容が最新である場合に、そのファイルの read 性能が向上することを確認した。今後は、他の分散ファイルシステムのキャッシュ機構について調査する。また、gfarm2fs のレイヤでのキャッシュ機構を設計・実装し、評価する予定である。

#### 謝辞

本研究を行うにあたって、有益な助言を頂いた筑波大学建部研究室の方々へ深く感謝する。また本研究は、科学技術振興機構戦略的創造研究推進事業 (JST CREST) の研究課題「ポストベタスケールデータインテンシブサイエンスのためのシステムソフトウェア」の支援を受けている。

#### 参考文献

- 1) Tatebe, O., Hiraga, K. and Soda, N.: Gfarm Grid File System, *New General Computing*, Vol.28, No.3, pp.257-275 (2010).
- 2) Szeredi, M.: FUSE: Filesystem in Userspace. <http://fuse.sourceforge.net/>.
- 3) 石黒 駿, 村上じゅん, 大山恵弘: Gfarm のためのカーネルドライバへのキャッシュ機構導入の検討, 並列/分散/協調処理に関するサマー・ワークショップ (SWoPP2011) (2011).