

# ext4 オンラインデフラグ実装に関する研究開発

藤田 朗                      佐藤 尚  
NECソフトウェア東北株式会社

## 概要

ext3 は Linux の標準ファイルシステムとして広く利用されてきたが、年々増加傾向にあるデータサイズに対応しきれなくなってきた。ext3 が抱える諸般の問題を解消する目的で、次世代ファイルシステムの ext4 が開発された。ext4 はできるだけ連続した領域のブロックを割り当てるよう考慮されているが、並列書き込み時や空き容量が少ない場合、フラグメントが発生しやすく、システム停止をせずにフラグメントを解消する方法がなかった。そこで我々は、オンラインでフラグメントを解消する方法を考案し、実装した。

本論文では、linux-2.6.31 の rc(リリース候補)版に含まれた ext4 オンラインデフラグ機能の実装とフラグメント解消による性能向上について述べる。

## 1. はじめに

ファイルシステムを長く運用している場合、それまでの書き込み、削除によりファイルデータのブロックがディスク上離れた領域に格納され、フラグメントが発生する。ファイルの読み込み時に複数箇所のブロックを検索するため、ハードディスクの無駄なシークが発生し、連続した領域に対するアクセスに比べ I/O 性能が悪化する。さらに、度重なるハードディスクのシークにより、駆動部分の耐久性が低下し、ひいてはハードディスク故障によるシステムの信頼性の低下にもつながる。

フラグメントが I/O 性能に及ぼす影響を図 1 に示す。フラグメントが発生した 1GB のファイルを読み込んだ場合、どのファイルシステム種別でも読み込み性能が劣化している。

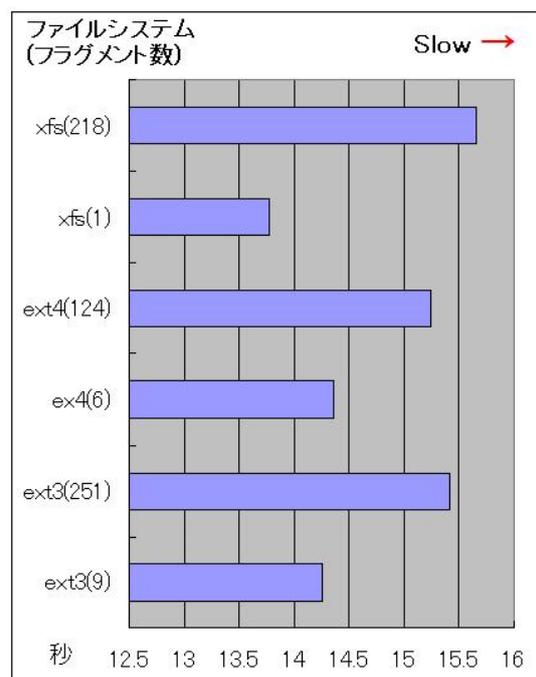


図 1 フラグメントによる読み込み性能劣化  
評価環境 カーネル: 2.6.31-rc3, アーキテクチャ: i386,  
CPU: Xeon 3.0GHz, メモリ: 1GB

先進的なファイルシステムでは、発生したフラグメントを解消し性能劣化を改善す

るためデフラグ機能に対応している(表1)。しかし、Linuxの標準ファイルシステムであるext3ではデフラグ機能に対応していない。つまり、ext3ではシステム停止をせずにフラグメントを解消する方法がなかった。

ext4がデフラグ機能に対応しない場合、Linuxユーザがext4採用を敬遠することも考えられる。さらに、デフラグ機能などの管理ツールが充実していないことは、WindowsからLinuxへの移行を検討する際の障害にもなりかねない。しかし、デフラグ機能を実装し、フラグメントを解消することでハードディスクのI/O性能が向上し、システム全体のスループットの改善も見込める。そこで我々は、ext4オンラインデフラグ機能を開発し、Linuxコミュニティへ提案した。

表1 主要ファイルシステム(FS)のデフラグ機能対応状況

FS 種別	デフラグ対応状況
ext3	×
ext4	○*
JFS	×
NTFS	○
XFS	○
Btrfs	○

○: 対応 ×: 未対応

上記は執筆時点の最新版であるlinux-2.6.31-rc3を参考にした。

※ ext4は我々が行った開発により、デフラグ機能に対応した。

## 2. ext4の研究開発の動向

ext4は機能強化と諸元拡大する目的として、Theodore Ts'o氏を中心に開発が進めら

れた。ext4では、最大ファイルサイズがext3の8倍になり(表2)、ブロックアロケーション機能の強化がなされている(表3)。

表2 ext3、ext4の諸元比較

FS 種別	最大ファイルサイズ	最大FSサイズ
ext3	2TB	16TB
ext4	16TB	1EB

表3 ext4の主要機能

機能	効果
Extents*	・ファイルサイズ、ファイルシステムサイズ拡大 ・ブロックの検索高速化
Delayed allocation	連続ブロック割り当て
Multi-block allocation*	連続ブロックでのブロック割り当て効率化
Flex_bg, Uninit_bg	ファイルシステム異常時の復旧時間の短縮化
Journal checksumming	ジャーナルの性能と信頼性の向上

※ ext4 オンラインデフラグに関する機能

### ●間接ブロック管理とエクステンツ管理

ext4で扱えるファイルサイズが大きくなったのに伴い、新しいファイル管理方式が導入された。

ext4ではext3で用いられてきた従来の間接ブロック管理(図2)のほか、標準機能としてエクステンツ形式でのファイルブロック管理(図3)をサポートしている。

間接ブロック管理では、0-11のブロックは直接ブロックであり、12,13,14番目のブロックはそれぞれ1段、2段、3段の間接ブ

ロックである。ファイルサイズが大きくなるにつれて、多くの間接ブロックを参照する必要があり、大きいファイルを扱うのには向いていない。

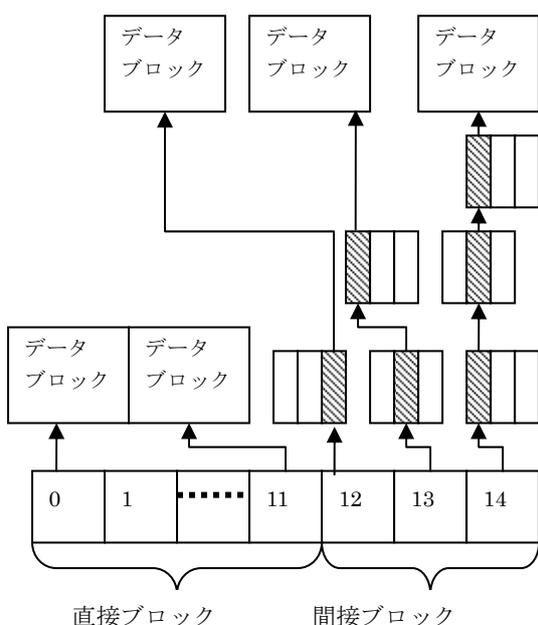


図2 間接ブロック管理

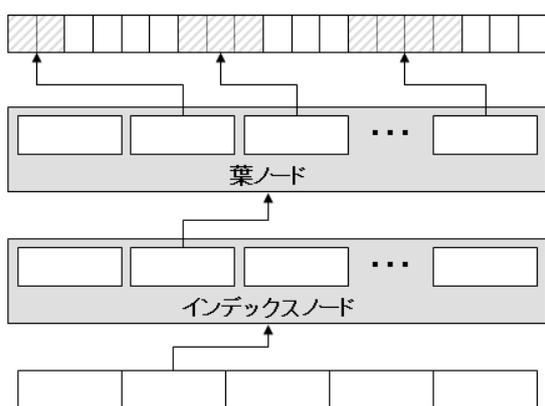


図3 エクステン形式での管理

エクステン形式では連続ブロックを領域として管理する。フラグメントが少ない場合、大きなファイルでもアクセス時間の低下は発生しにくい。ただし、フラグメン

トが悪化し、エクステン数が増加すると間接ブロック管理と同様に非効率なアクセスとなってしまふ。

表4はext3, ext4上に1GBのファイルを作成した際のフラグメント数による読み込み時間の性能劣化を示す。ext3と比べext4のフラグメント数は減少しているが、ブロック管理の方法に関係なく、フラグメントが読み込み速度に影響していることがわかる。

表4 フラグメントによる性能劣化

FS (管理方法)	フラグメント	読み込み(秒)	性能劣化
ext3(間接)	9	14.697	-8.1%
	251	15.894	
ext4(間接)	5	15.554	-10.9%
	123	17.246	
ext4(エクステン)	6	14.901	-3.4%
	125	15.404	

フラグメントの作成方法:dd コマンドを10並列で実行。

読み込み(秒)はcat コマンドを実行した5回分の平均時間。

### ●フラグメント

ext4ではサポートできる最大ファイルサイズの拡大に伴い、データ書き込み時に連続した領域にデータを格納できるようにするため、以下の2つの機能をサポートし、フラグメントの発生を抑制している。

- **Delayed allocation**  
ディスクへのデータの書き込みを遅延させ、ブロック割り当てのリクエストをまとめることで連続した領域にデータを書き込む。
- **Multi-block allocation**  
1度のブロック割り当て処理で連続した複数のブロックを割り当てる。

実際の運用時には、複数プロセスによる並列書き込み、空き容量の減少によりフラグメントが大量に発生しやすくなる。フラグメントによる I/O 性能の低下は、システム運用において軽視できるものではない。

ext4 は標準で大きなサイズを扱えるエクステンツ形式のファイルをサポートしている。従来の間接ブロック管理のファイルは互換性を維持するために利用可能となっているが、ファイルサイズの諸元も小さいことから、ext4 ではあまり利用されないことが見込まれる。そのため、ext4 オンラインデフラグ機能は、エクステンツ形式のファイルをサポートしている。

### 3. ext4 オンラインデフラグ実装方式

ext4 オンラインデフラグは新規に考案したユーザ空間プログラムの e4defrag コマンドとカーネル空間の EXT4\_IOC\_MOVE\_EXTENT ioctl により実現される機能である。

ext4 オンラインデフラグ機能は 2006 年後半から Linux コミュニティに提案活動を行ってきた<sup>1</sup>。約 3 年間におよびコミュニティで議論<sup>[3] [4] [5] [6]</sup>を繰り返した結果、2009 年 6 月にリリースした ver. 2.0 が正式に Linux に含まれることになった<sup>[7] [8]</sup>。

開発当初の ext4 オンラインデフラグは、デフラグに特化した ioctl を実装していた。コミュニティとの議論により、ioctl を汎用的にし、カーネル空間で行っていた donor

file(データ交換用ファイル)の作成、連続ブロックの割当てをユーザ空間で行うよう変更した。

ここでは、ext4 オンラインデフラグ ver. 2.0 の実装について述べる。

#### 3-1 e4defrag コマンドでのブロック交換

e4defrag コマンドによるユーザ空間からみたブロック交換処理を説明する。

- ① ユーザ空間で新規に donor\_file を作成する。作成した donor\_file に unlink システムコールを呼び出し、ファイルのリンク数を 0 にする。

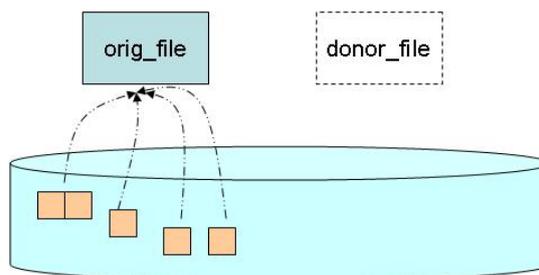


図 4 donor\_file の作成

- ② fallocate システムコール<sup>3</sup>を用いて donor\_file に対して、orig\_file(デフラグ対象ファイル)と同じ数の連続したブロックを割り当てる。

<sup>3</sup> fallocate システムコール

開始オフセットからバイト分を指定し、領域を対象ファイルに割り当てるシステムコール。割り当てられた領域にはデータを書き込まないため、高速に処理が終了する。デフラグではこのシステムコールを用い、あらかじめ連続したブロックをファイルに割り当て、後に呼ばれる EXT4\_IOC\_MOVE\_EXTENT でデータを書き込む。

<sup>1</sup> Linux ファイルシステム開発用のメーリングリストのアーカイブ、

<http://marc.info/?l=linux-fsdevel&r=1&w=2>

<sup>2</sup> ext4 開発用のメーリングリストのアーカイブ、

<http://marc.info/?l=linux-ext4&r=1&w=2>

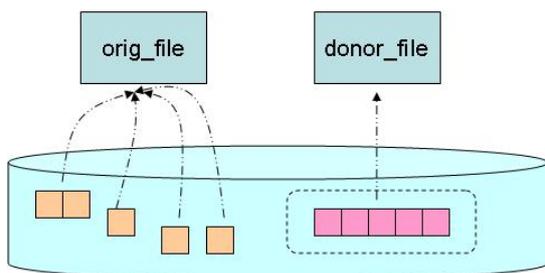


図5 fallocateによるブロック割り当て

- ③ FS\_IOC\_FIEMAP ioctl<sup>4</sup>を呼び出し、donor\_fileのエクステント数を算出する。②で割り当てたブロックのエクステント数(=フラグメント数)がorig\_fileより少ない、つまりデフラグによりフラグメントの改善が見込めることを確認する。

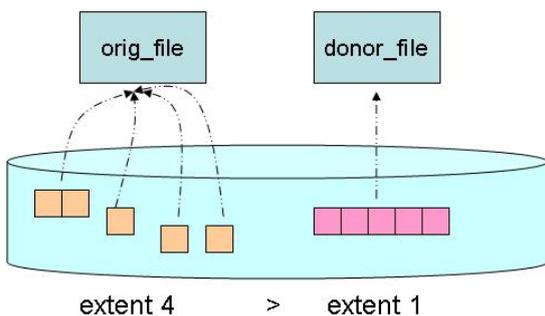


図6 エクステント数の比較

- ④ donor\_file のエクステント 1 個ごとに EXT4\_IOC\_MOVE\_EXTENT ioctl を呼び、orig\_file、donor\_file 間でデータブロックを交換する。この結果、orig\_file と donor\_file が指す物理ブロックが交換され、orig\_file はフラグメン

<sup>4</sup> FS\_IOC\_FIEMAP ioctl

指定したファイルの論理オフセット、物理オフセット、ブロック長などのエクステント情報を取得する ioctl。取得したエクステント数はフラグメント数に等しいため、e4defrag コマンドでは orig\_file, donor\_file のフラグメント数を比較し、改善が見込める場合のみ処理を継続する。

トが少ない連続したブロックを保持することになる。

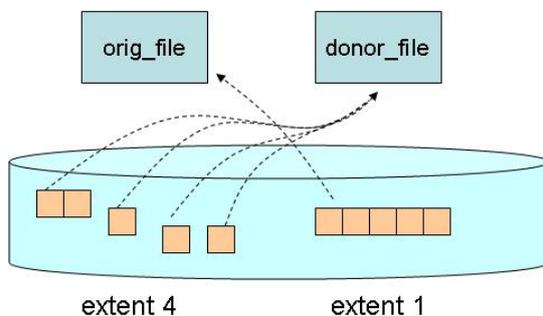


図7 ブロック交換

- ⑤ ブロック交換後の donor\_file は必要ないため、後処理として close システムコールを呼び出す。リンク数が 0 である donor\_file は削除される。

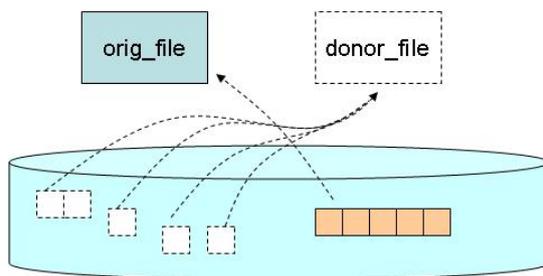


図8 donor\_file の削除

### 3-2 EXT4\_IOC\_MOVE\_EXTENT

EXT4\_IOC\_MOVE\_EXTENT は指定した 2 つのファイル間でブロックを交換する ioctl である。move\_extent 構造体(図 9)を入出力引数とし、ファイルの論理オフセット orig\_start から len ブロック分をデフラグ対象ファイル(orig\_file)とデータ交換用の一時ファイル(donor\_file)間で交換する。

```

#define EXT4_IOC_MOVE_EXT      _IOWR('f', 15, struct move_extent)

struct move_extent {
    __u32 reserved;           /* should be zero */
    __u32 donor_fd;          /* donor file descriptor */
    __u64 orig_start;        /* logical start offset in block for orig */
    __u64 donor_start;       /* logical start offset in block for donor */
    __u64 len;               /* block length to be moved */
    __u64 moved_len;        /* moved block length */
};

```

図9 EXT4\_IOC\_MOVE\_EXTENT と move\_extent 構造体

EXT4\_IOC\_MOVE\_EXTENT ioctl のカーネル内部で行うブロック交換処理について説明する。

- ① move\_extent 構造体に設定された値を元に、orig inode のデータを交換する論理オフセットの範囲を決定する。

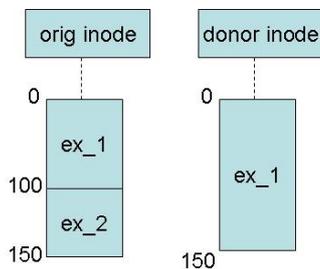


図10 ブロック交換範囲の設定

※ 上記の例では orig inode の先頭ブロックから 150 ブロック目までの 151 ブロックを対象とする。

- ② 指定された論理オフセットに該当する orig inode のファイルデータをページキャッシュに読み込む。

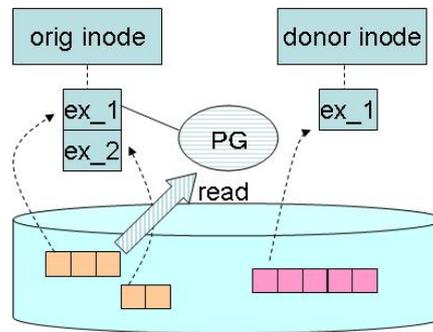


図11 ページキャッシュ読み込み

- ③ データブロック交換用に orig inode, donor inode とともにエクステントを複製する。
- ④ ブロック交換処理はページサイズごとに行うため、複製したエクステントを 1 ページサイズ分に整形する。

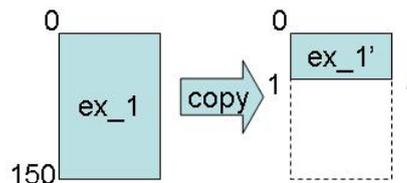


図12 エクステントの複製とサイズ調整

- ⑤ 整形したエクステントを orig inode の

同じ論理オフセットに該当するエクステンツトへ上書きする。これにより、**orig** のエクステンツトが指すブロックが変更される。

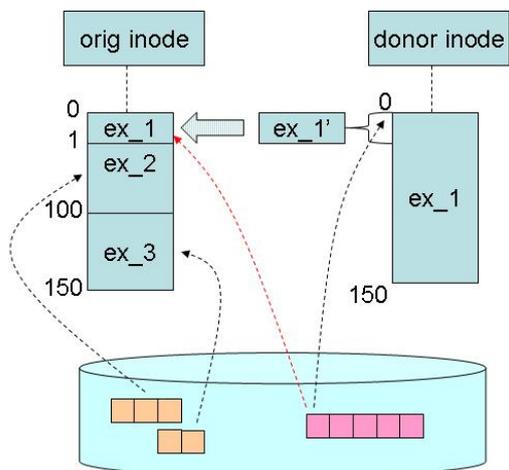


図 1 3 ブロック交換

- ⑥ ④、⑤で行った処理を今度は **orig inode** から **donor inode** へ実行する。これによりお互いのブロックがページサイズ分交換される。
- ⑦ ページキャッシュ上の変更されたデータをディスクへ書き出す。**orig inode** のデータが **donor inode** のブロックへ書き出される。

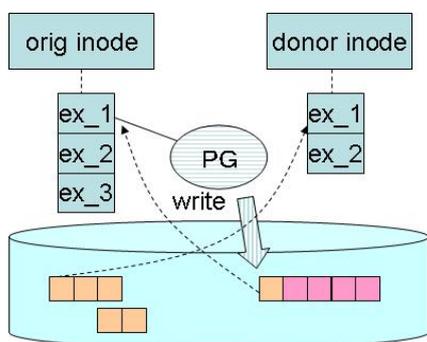


図 1 4 ページキャッシュ書出し

②-⑤の処理を 1 ページごとに実施し、①で算出した範囲(例では先頭から 150 ブロック目まで)を実施するまで繰り返す。

### 3-3 ファイルアクセスの排他

オンラインでデフラグを実施するために、カーネル空間では **inode** 構造体の **mutex** ロック (**i\_mutex**) と **ext4\_inode\_info** 構造体のセマフォ (**i\_data\_sem**) を利用する。

**i\_mutex** は他プロセスからの **truncate** を防ぐために取得され、カーネル空間での処理が終了するまで保持する。

**i\_data\_sem** はブロックを参照する処理では読み込みセマフォ、データブロックの割り当て・削除を行う処理では書き込みセマフォを局所的に取得・解放する。

XFS のデフラグコマンドの **xfs\_fsr** はデフラグ実行中にファイルが変更された場合は処理が失敗となってしまいが、**ext4** オンラインデフラグではこれらを用いた排他により、他のプロセスが占有しているファイルに対してもデフラグが実行可能となっている。

### 3-4 性能改善結果

**ext4** 上に作成した **2GB** のファイルに対し、**ext4** オンラインデフラグ実行前後のフラグメント数と読み込み時間の改善具合を表 5 に示す。デフラグの実行により読み込み時間が約 **16%** 改善している。

表5 デフラグによる性能改善

デフラグ	フラグメント	読み込み(秒)	改善
前	257	39.893	15.6%
後	24	33.653	

評価環境 カーネル: 2.6.31-rc3, アーキテクチャ: i386,

CPU: Xeon 3.0GHz, メモリ: 1GB

性能測定方法: dd コマンドを 10 並列で実行し作成した 2GB のファイルに対し、ext4 オンラインデフラグ実行前後での cat コマンドの読み込み時間を測定。

表 5 は 5 回測定した平均値。

#### 4. 今後の課題

ext4 オンラインデフラグ機能を実装することにより、ファイル単位のデフラグが可能となった。今後は、アプリケーションの性能向上、連続空きブロックが少ない状態でのデフラグを可能にするため、以下の 2 つの機能の追加を計画している。

1. 指定したディレクトリ配下のファイルを物理的に近い領域に配置する
2. ファイルシステムの空きが少ない状態でも指定したファイルのフラグメントを優先的に解消する

1 はアプリケーションが同時にアクセスするファイルを近隣に配置し、アプリケーションの性能向上を図る機能である。

2 は連続した空きブロックが少ない状態でも、デフラグ非対象ファイルを別の領域に移動することで連続した空き領域を作成し、そこにデフラグ対象ファイルを格納する機能である。これにより、ext4 上に十分な連続空きブロックがない状態でも、アク

セス頻度が高いファイルなどに対し、優先的にデフラグを実行することが可能となる。

これら 2 つの機能を ext4 オンラインデフラグに追加するためには、既存の ext4 のブロックアロケーションをより柔軟で高機能なものにしなければならない。そのため、我々は Linux コミュニティへブロックアロケーション制御機能の提案を実施した<sup>9</sup>。

ブロックアロケーション制御機能は、ext4 上の指定した領域への優先的なブロック割り当て、または特定の領域をブロックアロケートから保護することを可能にする。現在、ファーストリリースに対して得られた Linux コミュニティからの意見を反映し、次のリリースに向けての開発を行っている。

#### 5. まとめ

ext4 ファイルシステムのフラグメントを解消する方法として ext4 オンラインデフラグ機能を考案し、実装評価した結果、フラグメント解消による性能向上を実現した。

継続的な提案活動と、Linux コミュニティの協力もあり、ext4 オンラインデフラグ機能は 2009 年 6 月に、linux-2.6.31-rc1 にマージされた<sup>5</sup>。今後問題がなければ次のバージョンで Linux の本流へ取りこまれる予定である。また、ユーザ空間プログラムの e4defrag コマンドも 2009 年 7 月に e2fsprogs-1.41.8 に含まれた<sup>6</sup>。

<sup>5</sup> Linuxに取り込まれたパッチのcommit log,

<http://git.kernel.org/?p=linux/kernel/git/torvalds/linux-2.6.git;a=commit;h=748de6736c1e482e111f9d1b5a5d5b1787600cad>

<sup>6</sup> e2fsprogsに取り込まれたパッチのcommit log,

<http://e2fsprogs.git.sourceforge.net/git/gitweb.cgi?p=e2fsprogs;a=commit;h=4836459a3fe8ff208a7bf21e28e4648bdde76440>

カーネル、コマンドともに正式なバージョンに含まれたことにより、今後多くの利用者からのフィードバックが得られ、より品質やユーザビリティの強化が図られることが期待される。

今後も ext4 への新機能の提案や開発、品質向上の活動を通じて、Linux コミュニティに貢献していきたい。

## 参考文献

[1] Takashi Sato. ext4 online defragmentation.

*In Ottawa Linux Symposium (2007)*

<https://ols2006.108.redhat.com/2007/Reprints/sato-Reprint.pdf>

[2] Daniel P. Bovet, Marco Cesati.

詳解 LINUX カーネル 第3版 (2007),

ISBN 978-4873113135

[3] Akira Fujita.

“[RFC][PATCH 0/3] ext4: online defrag (ver 1.0),”

<http://marc.info/?l=linux-ext4&m=123329594919001&w=2>

[4] Chris Mason.

“Re: [RFC][PATCH 0/3] ext4: online defrag (ver 1.0),”

<http://marc.info/?l=linux-ext4&m=123334655401020&w=2>

[5] Greg Freemyer.

“Re: [RFC][PATCH 0/3] ext4: online defrag (ver 1.0),”

<http://marc.info/?l=linux-ext4&m=123335484515745&w=2>

[6] Theodore Ts'o.

“Re: [RFC][PATCH 1/3] ext4: Add EXT4\_IOC\_DEFRAG ioctl and basic defrag functions,”

<http://marc.info/?l=linux-ext4&m=123375929109850&w=2>

[7] Akira Fujita.

“[RFC][PATCH 0/3]ext4: online defrag (ver 2.0),”

<http://marc.info/?l=linux-ext4&m=124297598014251&w=2>

[8] Theodore Ts'o.

“Re: [RFC][PATCH 1/3] Add EXT4\_IOC\_MOVE\_EXT ioctl and related functions,”

<http://marc.info/?l=linux-ext4&m=124489932017806&w=2>

[9] Akira Fujita.

“[RFC][PATCH 0/7]ext4: Block allocation restriction and inode preferred range of blocks for ext4,”

<http://www.spinics.net/lists/linux-ext4/msg13790.html>