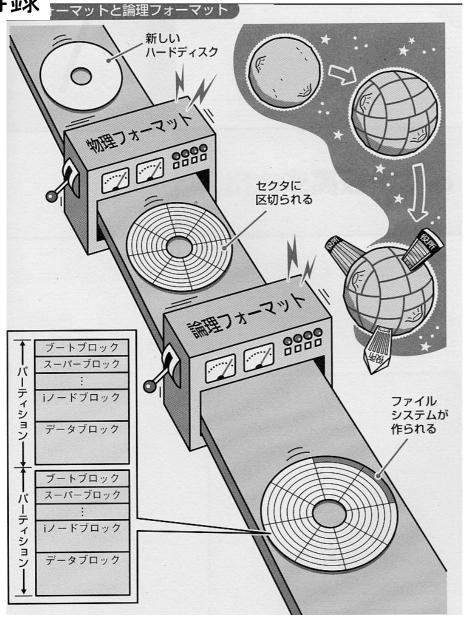
オペレーティングシステム2004 ファイル管理(2)

2003年11月11日 海谷 治彦

目次

- ・記録装置の概要(復習)
- FAT
 - ― 旧Windowsのファイルシステムだが、USBメモリ等では未だにコレが使われることがある。
 - こっちのほうが簡単だから先に話します.
- Ext2
 - Linuxで最も一般的なファイルシステム.
- Ext3
 - 昨今ではこっちのほうが一般的かも
 - ジャーナルファイルシステム

再録



フォーマット

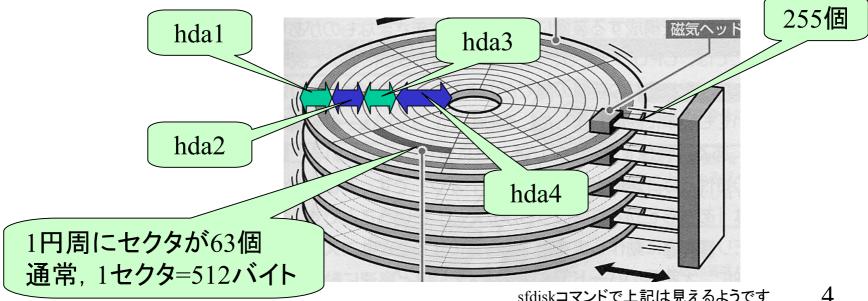
- 通常、隣接したシリンダ何枚かをグループ化し、それをパーティションとする。
- ディスク丸ごと1パー ティションとしてもさ しつかえない。

再録 とあるディスクの情報の実例

ディスク /dev/hda: ヘッド 255, セクタ 63, シリンダ 2434 ユニット = シリンダ数 of 16065 * 512 バイト デバイス 始点 終点 ブロック ID システム /dev/hda1 1 255 2048256 83 T_iinux

4つのパ-ティション

82 Linux スワップ 530145 /dev/hda2 256 321 /dev/hda3 322 576 2048287+ 83 Linux /dev/hda4 577 2434 14924385 83 Linux



データ読書きの補足

- データは基本的にセクタ単位(512Bもしくは 1024B単位)で物理的に読書きする.
- 大抵, どの論理フォーマット(ファイルシステム)でも、複数のセクタをグループ化して扱う.
 - FATの場合, クラスタ
 - Ext2の場合, ブロック と呼んだりする.
- 理由はセクタが小さすぎるからだろう。

FAT (File Allocation Table)

- 論理フォーマットの一種
- 旧Windowsで利用されていた.
- 今でも利用されているところがある.
- Linuxでも読書きできる.



左記は私が使ってる USBメモリ

FATの内部形式

- 連続した2ⁿ個のセクタをグループ化してクラスタとする.
 - とはいえディスクの先頭部分は例外
- ファイルは1個以上のクラスタに分割して配置されている。
- ファイルがどこにあるかの情報はFATとディレクトリエントリという情報で管理されている.

FATの内部構造

あるパーティション(論理ディスク)

MBR ブートセクタ **FAT** 予備 ルートディレクトリ エントリ クラスタ1 クラスタ2 クラスタ3 クラスタn

この辺はファイルの情報とは関係無し

FATテーブル

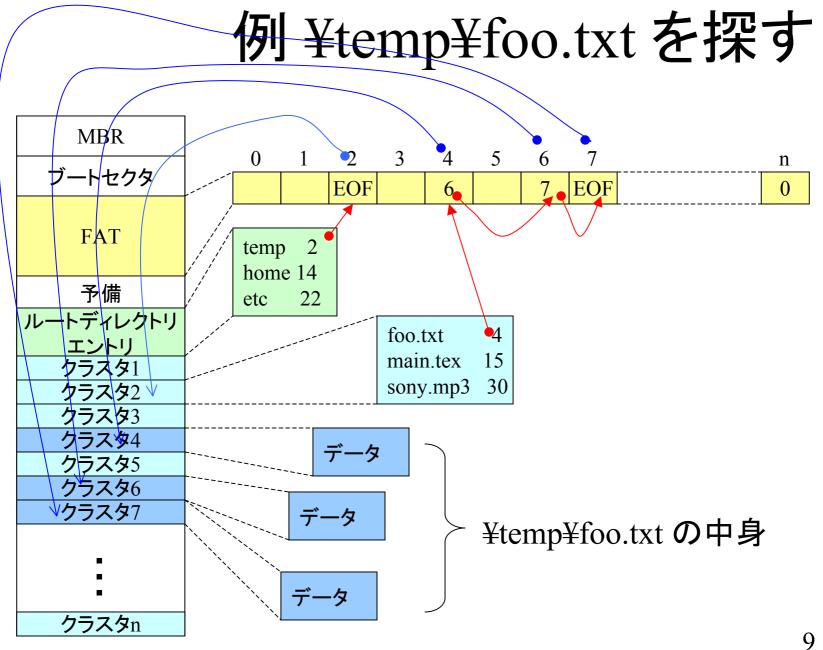
下にあるクラスタの中身どのようにグルーピングされているかを管理.値としては、

- •未使用
- •EOF
- 続くクラスタの番号 のどれか。

ここのクラスタ毎に、

ファイルの中身や、

ディレクトリエントリ(サブディレクトリ内のファイルの情報) なんかが入っている.



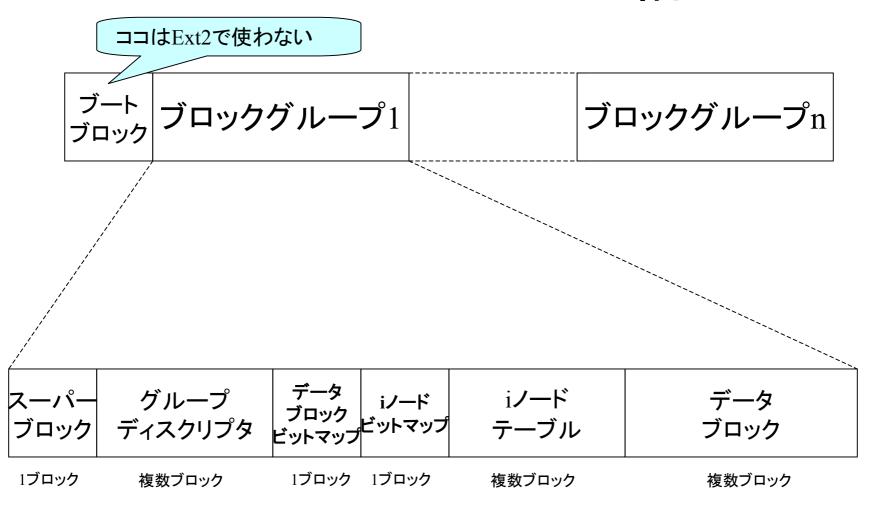
FATの特徴

- 仕組みが単純(?)
- (後述のext2でもそうだが)クラスタサイズが 効率に影響する.
 - 小さいと大きなファイルを扱うのが不便.
 - ・沢山のクラスタに実データが分散するため。
 - 大きいと小さいファイルを扱う場合無駄になる.
- ファイルの中身が複数のクラスタに分散し、 クラスタが近接していない場合、アクセス が遅くなる。
 - いわゆる fragmentation (断片化)というヤツ.

Ext2 拡張ファイルシステム2

- 論理フォーマットの種類の1つ
- HDだけでなく、例えばUSBメモリ等もExt2でフォーマットできる.
- Linuxでは最も標準的なファイルシステムの種類.
- パーティションをブロックという単位で分割して管理する.
 - 通常, 1ブロック 1024B~4096B
 - 1セクタが512Bくらいなので、2~8セクタをグループ化
- ブロックをグルーピングし、同一グループのブロックは隣接トラックに配置される。よって、同一ブロック内のデータは短い時間間隔でアクセスできる。
 - ブロックグループ

パーティション内の構造



ブロックグループ内の項目

- スーパーブロック (1ブロック)
 - パーティションの情報が書いてある.
 - 各ブロックグループに複製を持っており、実際使うのは、グループ0のもののみ。
- グループディスクリプタ (複数ブロック)
 - グループの情報が書いてある.
 - 各ブロックグループに複製を持っており、実際使うのは、グループ0のもののみ。
- データブロックビットマップ (1ブロック)
 - データブロックが使用されているか否かを0/1情報で記述したビット列.
- iノードビットマップ (1ブロック)
 - データブロックと同様.
- iノードテーブル (複数ブロック)
 - 後述.
- データブロック (複数ブロック)
 - 個々のファイルの中身が入っているブロック郡

ビットマップとサイズ制限

- ビットマップは1ブロックに収めなければならないので、結果としてブロックやi nodeの数に制約がある.
- 一般に、ブロックサイズは1024B~4096B、すなわち、 8192bit ~ 32768bit
- よって、ブロックサイズを1024Bにした場合、
 - グループ内のブロック数は8192個
 - グループ内のi node数上限も8192個となる。
- ま、ユーザーからはグループは認知できないので、作成してよいファイル数の上限とかが気になることはないだろう。
 - とはいえ、1つのパーティションに作成できるファイル数には上限 はあるよ。

iノード

- Ext2ファイルシステム内のファイルはiノード番号で区別されている。
 - ファイル名は相対的な名前に過ぎない.
- ファイル固有の情報は、
 - ファイルの種類とアクセス権 (16bit)
 - 所有者のID (16bit)
 - バイト数 (32bit)
 - ハードリンクのカウンタ (16bit)等がある.
- 個々のiノードの情報は128バイトの構造体にまと められている.

```
ext2 inodeの一部
struct ext2 inode {
    u16 i mode;
                   /* File mode */
     u16 i uid;
                  /* Owner Uid */
    _u32 i_size;
                  /* Size in bytes */
    u32 i atime; /* Access time */
     u32 i ctime; /* Creation time */
    __u32 i_mtime; /* Modification time */
    u32 i dtime; /* Deletion Time */
    u16 i gid; /* Group Id */
    u16 i links count; /* Links count */
    _u32 i_blocks; /* Blocks count */
    u32 i flags; /* File flags */
   union {
       // ** 省略
   } osd1;
                  /* OS dependent 1 */
     u32 i block[EXT2 N BLOCKS];/* Pointers to blocks */
     u32 i version; /* File version (for NFS) */
     _u32 i_file_acl; /* File ACL */
    u32 i dir acl; /* Directory ACL */
                                                 include/linux/ext2 fs.h
     u32 i faddr; /* Fragment address */
   union {
                                                 の217行目あたりから
       // ** 省略
                /* OS dependent 2 */
全部で128バイト
    } osd2;
};
```

16

構造体 ext2 inode の概要

- ・ 個々のinode情報を示す構造体
- この構造体のインスタンスが、inodeテーブルのブロックに詰まっている。
- 1つのinodeインスタンスは128バイト.
- よって、1ブロックが1024Bなら、ブロック当たり8個のinodeのインスタンスが入っている。(意外に少ないね。)

iノードテーブル

- 前述のiノードを示す構造体のインスタンス (1個128バイト)が並んでいるテーブル.
- ・ ブロックサイズが1024Bなら,
 - 1ブロックに1024/128=8個のインスタンスが入る.
 - iノードビットマップは8192個のbitを持てる.
 - 別に最大bit数を利用しなくても良いが、利用したとすると・・・
 - iノードテーブルのために、8192/8=1024ブロックが必要となる。

Ext2でのファイルへたどり着く手順

- ファイル名 (パス名)から、そのファイル (が指す実体)のiノード番号を得る.
- 2. iノード番号をもとに, iノードテーブル内の 該当する構造体ext2_inodeのインスタン スを得る.
- 3. 構造体メンバーにファイル(の中身)が格納されているデータブロックの番号配列 (__u32 i_block[EXT2_N_BLOCKS])があるので、それに従い、ファイルの中身をデータブロックから引きずり出す.

ext2 inodeの一部

```
struct ext2 inode {
    u16 i mode;
                     /* File mode */
     u16 i uid;
                     /* Owner Uid */
     u32 i size;
                    /* Size in bytes */
                    /* Access time */
     u32 i atime;
     u32 i ctime;
                    /* Creation time */
    u32 i mtime; /* Modification time */
    u32 i dtime; /* Deletion Time */
    u16 i gid;
                     /* Group Id */
     u16 i links count; /* Links count */
    u32 i blocks; /* Blocks count */
                    /* File flags */
    u32 i flags;
    union {
        // ** 省略
                       /* OS rependent 1 */
    } osd1;
     _u32 i_block[EXT2_N_BLOCKS];/* Pointers to blocks */
     u32 i version; /* File version (for NFS) */
     _u32 i_file_acl; /* File ACL */
     _u32 i_dir_acl; /* Directory ACL */
     u32 i faddr; /* Fragment address */
    union {
        // ** 省略
                       /* OS dependent 2 */
    } osd2;
                 全部で128バイト
};
```

ここが実際のデータが入って いるブロックを格納している 配列. EXT2_N_BLOCKS は 15 コードの181行あたり. 後述のように最後の3つが間 接,二重間接,三重間接ブロッ クに使われている.

include/linux/ext2_fs.h の217行目あたりから

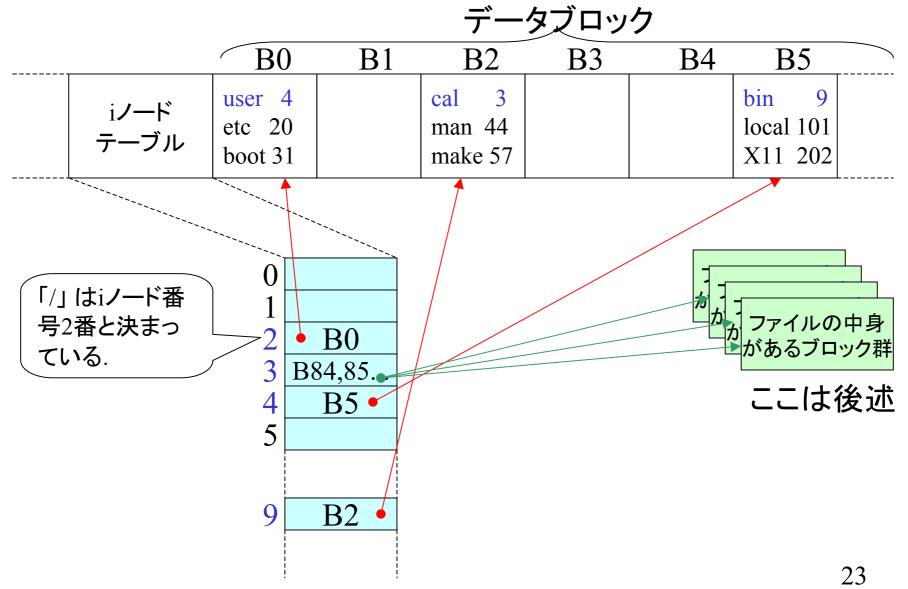
iノードとファイルの種類

- Linux(UNIX)では全てのファイルにiノード (番号)がある。
- ディレクトリもシンボリックリンクもiノード番号を持つ。
- ことを思い出しておいてください.

パス名からiノード番号を得る

- 1. ルードディレクトリ/のiノード番号は「2」と決まっているので、現在の注目するiノードを2に設定する.
- 2. 注目しているiノードの構造体インスタンスに指定されているデータブロックの中身を見て、ディレクトリ内にあるファイル名とiノード番号の対応表を得る.
 - a. 知りたいパスがディレクトリ内にあれば,対応する番 号が結果.
 - b. まだパスの途中なら,パスの途中であるディレクトリに対応するiノード番号を注目するiノードに設定しなおし,2に戻る.

例: /usr/bin/cal の番号を探す



iノード構造体とファイルの中身

- 前述のようにiノード構造体は(たった)128B しかない。
- iノードが指しているファイルやディレクトリの中身は別途、データブロックに格納されている。
 - ディレクトリの場合、属しているファイルのリスト
- 中身が入っているデータブロックのありかは構造体のメンバ, i_block[] 配列に記録されている。

ext2 inodeの一部

```
struct ext2 inode {
    u16 i mode;
                     /* File mode */
     _u16 i_uid;
                     /* Owner Uid */
     u32 i size;
                    /* Size in bytes */
                    /* Access time */
     u32 i atime;
     u32 i ctime;
                    /* Creation time */
    u32 i mtime; /* Modification time */
    u32 i dtime; /* Deletion Time */
    u16 i gid;
                     /* Group Id */
     u16 i links count; /* Links count */
     u32 i blocks; /* Blocks count */
                    /* File flags */
    u32 i flags;
    union {
        // ** 省略
                       /* OS rependent 1 */
    } osd1;
     _u32 i_block[EXT2_N_BLOCKS];/* Pointers to blocks */
     u32 i version; /* File version (for NFS) */
     _u32 i_file_acl; /* File ACL */
     _u32 i_dir_acl; /* Directory ACL */
     u32 i faddr; /* Fragment address */
    union {
        // ** 省略
                       /* OS dependent 2 */
    } osd2;
                 全部で128バイト
};
```

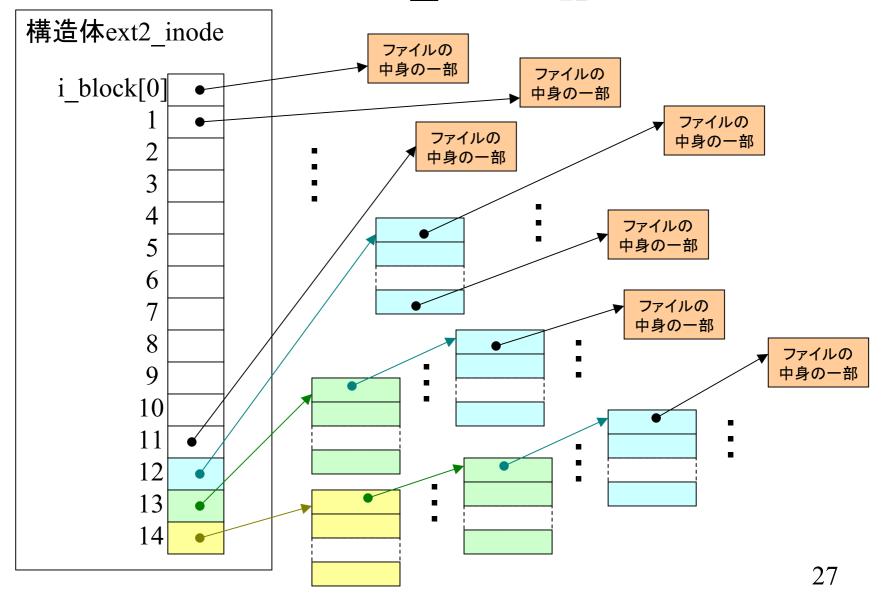
ここが実際のデータが入って いるブロックを格納している 配列. EXT2_N_BLOCKS は 15 コードの181行あたり. 後述のように最後の3つが間 接,二重間接,三重間接ブロッ クに使われている.

include/linux/ext2_fs.h の217行目あたりから

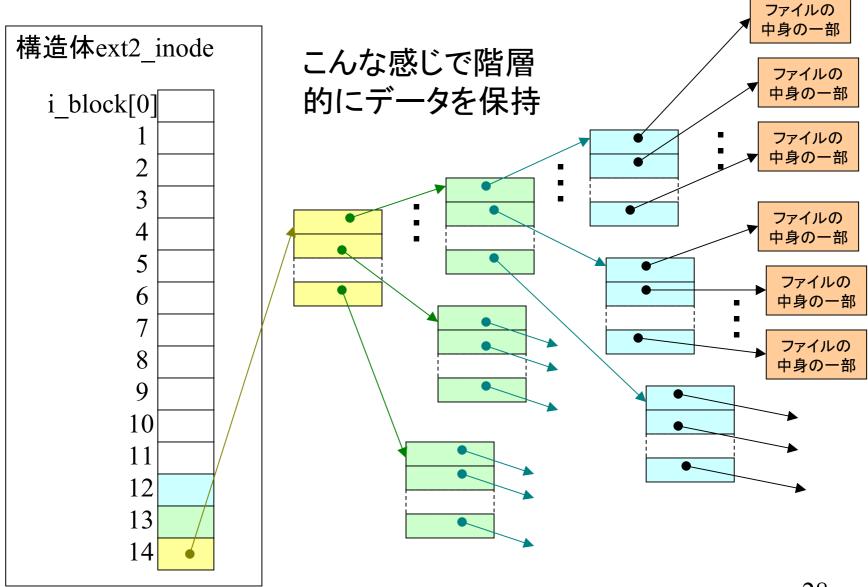
メンバi block[]

- EXT2_N_BLOCKS個の配列, 通常この値 は15.
- ・配列要素は以下の4種類に分かれる.
 - i_block[0]~[11]の12個: 直接アドレッシング
 - ・ファイルの中身が入るデータブロックを直接指す.
 - i_block[12] 一段間接アドレッシング
 - i block[13] 二段間接アドレッシング
 - i block[14] 三段間接アドレッシング
 - 間接アドレッシングについては後述

iノード構造体とi_block[]のイメージ



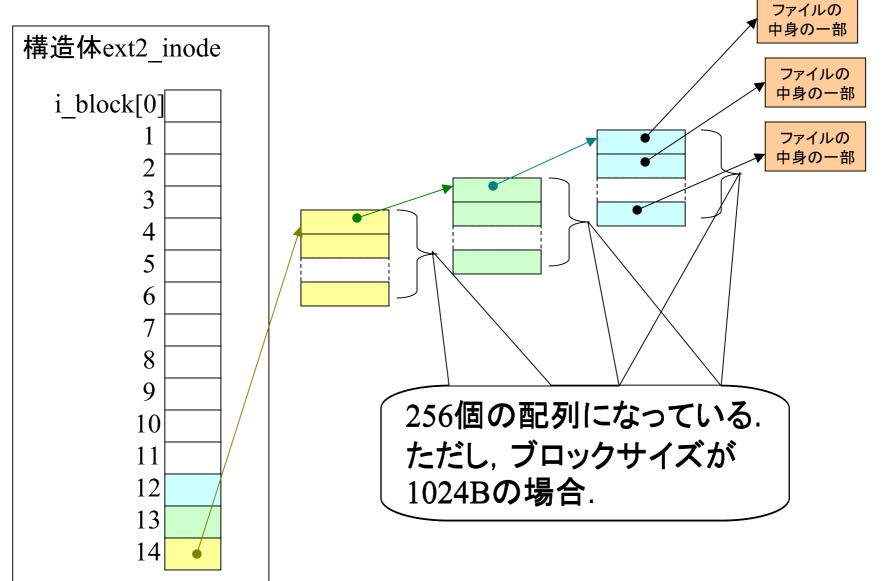
三段間接アドレッシング



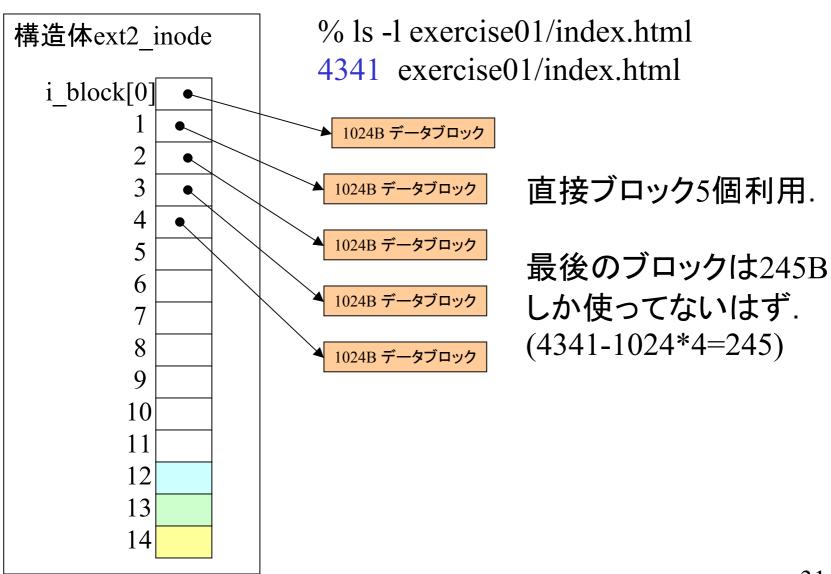
扱えるファイルサイズは?

- 1ブロックを、1024Bとすると、間接ブロックを使わないと、1024*12≒12KBのファイルしか扱えない.
- 間接ブロック内に1つのブロックの位置を保存するには4B必要らしいので、間接ブロックからは、1024/4=256個のブロックを指すことができる.
- ・ 13個目以降の要素では,
 - 間接 256*1024B≒262KB
 - 二重 256*256*1024B≒67MB
 - 三重 256*256*256*1024B≒17GB
 - 結構大きなファイルが扱えますな.
- 実際はハードウェアの制約等により、際限なく大きなファイルが扱えるわけでない。

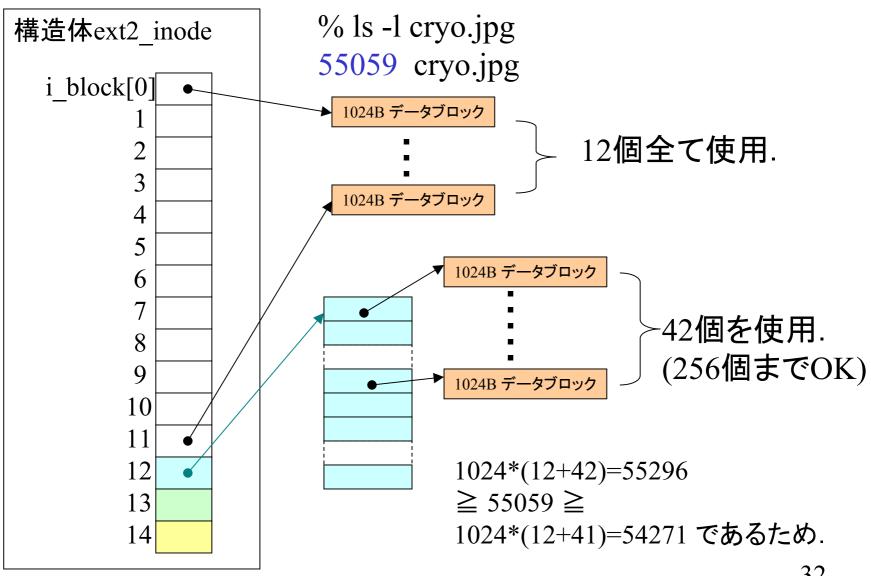
三段間接アドレッシング



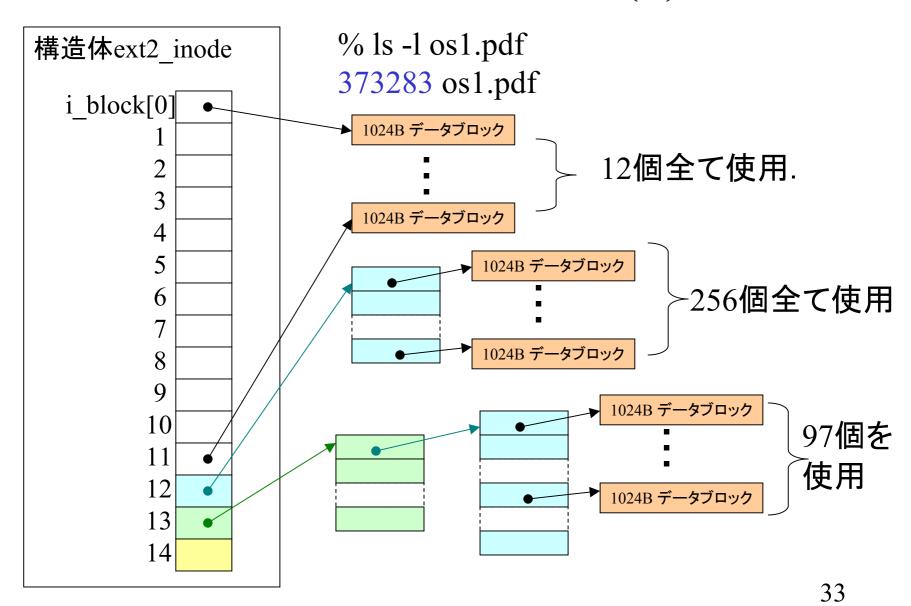
ファイルの格納例 (1)



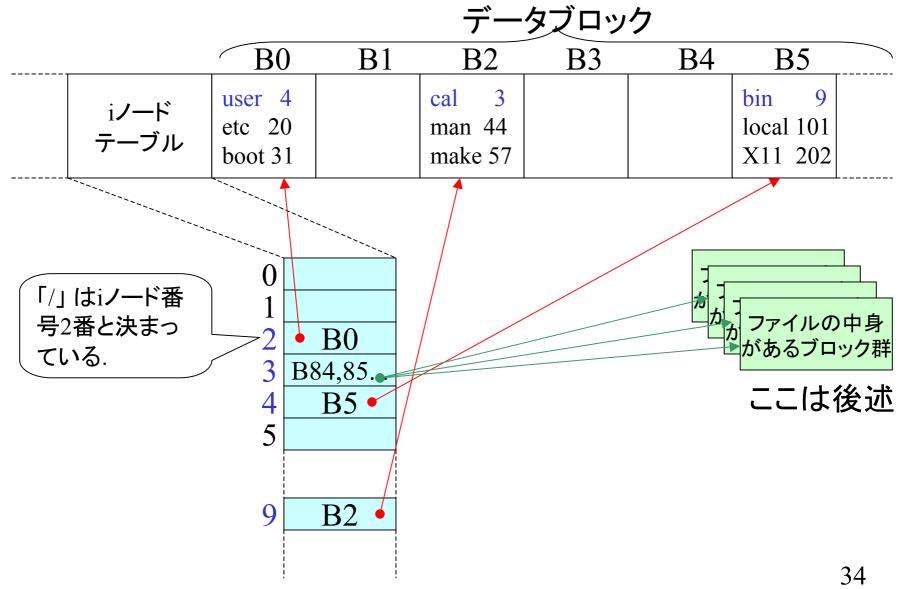
ファイルの格納例(2)



ファイルの格納例 (1)



ディレクトリの中身もデータブロック



ディレクトリのデータブロック

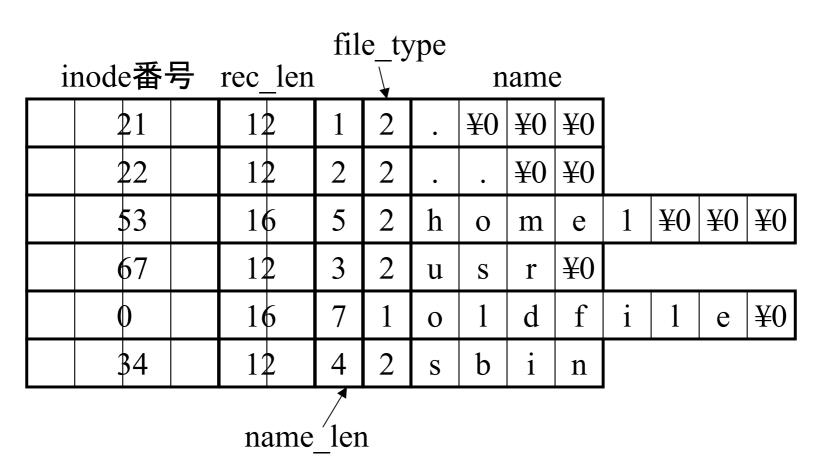
- 構造体 ext2_dir_entry_2 のインスタンスが 詰まっている.
- 上記構造体の一つ一つが、そのディレクト リ内にある他のファイルの情報を保持している。
- この構造体は可変長.
 - 最後のメンバーが文字列であるため.

ext2_dir_entry_2

型: 1 通常ファイル, 2 ディレクトリ, 7 シンボリックリンク 他

効率化のため4の倍数長になっている. 不要な 部分には ¥0 文字が詰めてある.

例: とあるデータブロックの中身



シンボリックリング

- リンク先のパスが60バイト以下の場合, ext2_inodeのメンバーであるi_block[]に埋め込む。
 - i_block[]は1個32ビット(4バイト)であるため, 15×4=60個.
- 60バイトを超える場合は、普通のファイル 同様、データブロック内に保存する.

Ext3

- ・ Ext2と互換性のあるジャーナリングファイルシステム
- 昨今のLinuxでは標準らしい.

ジャーナリング・ファイルシステム

- ・時間のかかるファイルシステムの整合性 チェックを短時間に行うための仕組み.
- 具体的には最近の更新情報をジャーナルと呼ばれる場所に格納し、整合性チェックを高速化する。

ファイルシステムの整合性

- ブロックの内容は通常、メモリに複製を作り(キャッシュと呼ぶ)、そちらを操作することで、高速化をしている。
- 実際のブロックの内容とキャッシュは最終的(システム停止時等)には一致していないといけない。
- これが一致しているかどうかをチェックする のが整合性.

Ext3の戦略

- ブロック単位で処理を行う.
- ジャーナルという特別なファイル領域を準備しておく。
- キャッシュのディスクへの書き戻しを以下 の三段階で行う。
 - 1. <u>ジャーナルへのコミット</u>: キャシュをまずジャーナルに書き込む.
 - 2. <u>ファイルシステムへのコミット</u>: シャッシュを実際のファイルシステムに書き込む.
 - 3. ジャーナルに書いたものを破棄する.

何故, 前述の戦略が良いか?

- ジャーナルへのコミット中にクラッシュが起きた場合
 - キャシュの更新自体、無かったことにする.
 - すなわち、ここ最近のファイル更新はパーになる.
 - しかし、ファイルシステムの一貫性は保たれる.
- ファイルシステムへのコミット中にクラッシュした 場合
 - ジャーナルに更新結果があるのでそれをコピーすればよい。
 - よって更新結果を完全に復旧できる.
- 要はファイルシステム本体に中途半端な更新情報が書き込まれるのを防ぐことができる。

どこまでジャーナルに入れるか?

以下の三種類があるらしい.

- journalモード
 - 全てのブロックをジャーナルに一度入れる.
 - 無論, 遅いが安全.
- ordered-−F
 - データブロック以外(メタデータと呼ばれる)のブロック のみジャーナルに入れる.
 - ファイルシステムの構造(内容ではない)の安全性を重視.
- writeback ₹—ド
 - ファイルシステムが更新したメタデータのみをジャーナルに入れる。