# オペレーティングシステム2004 ファイル管理(1)

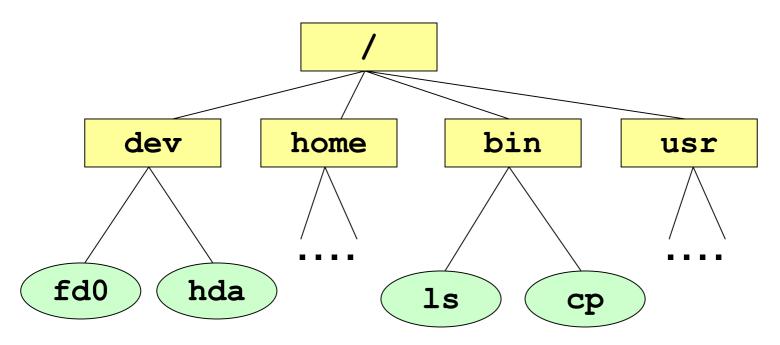
2003年11月5日 海谷 治彦

#### 目次

- (UNIX系)ファイルシステムの概要
- プロセスとファイルの関係
- 仮想ファイルシステム VFS
  - C言語 関数へのポインタ
- ディスクの物理構造とフォーマット

### UNIX系ファイルシステムの概要

ご存知のとおり、UNIX系OS(その他も大抵そうだけど)は、階層的なファイルシステムを持っている.

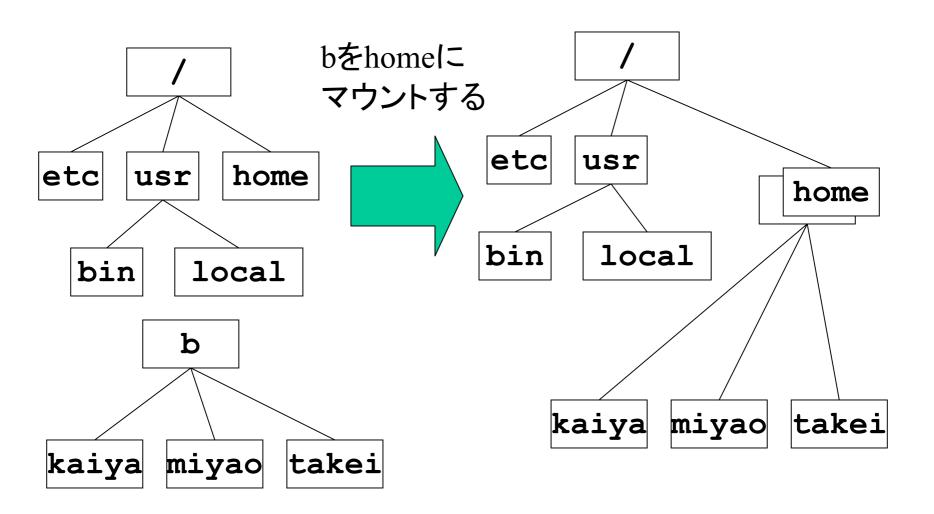


文献5 p.15

#### mount: ファイルシステムの接木

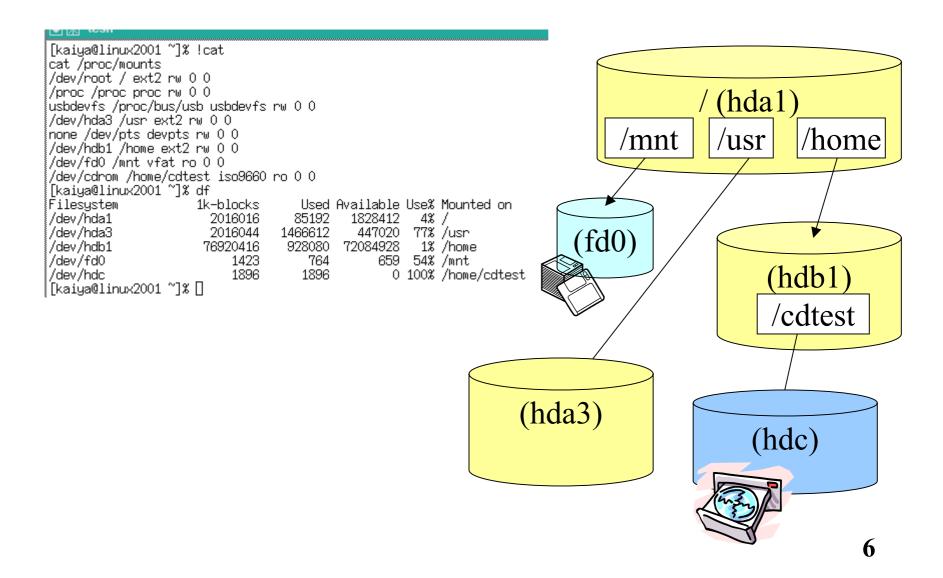
- 実際のファイルシステムは、複数のファイルシステムを接木して、一つの大きな木構造を成している。
- これによって、物理的/論理的に異なるファイルシステム群を1つの構造にまとめることができる。
  - 物理的: HDやFDやCDやネットワークデバイス
  - 論理的: Ext2やFATやNTFSやSMB等
    - Ext2はLinux標準のファイルシステム
    - FATは古いWindows(MSDOS)のファイルシステム
    - NTFSはWin/NtやXPのファイルシステム
    - SMBはWindowsのネットワークドライブ

## mountの概念図



文献1 p.45 改

### mountの実際



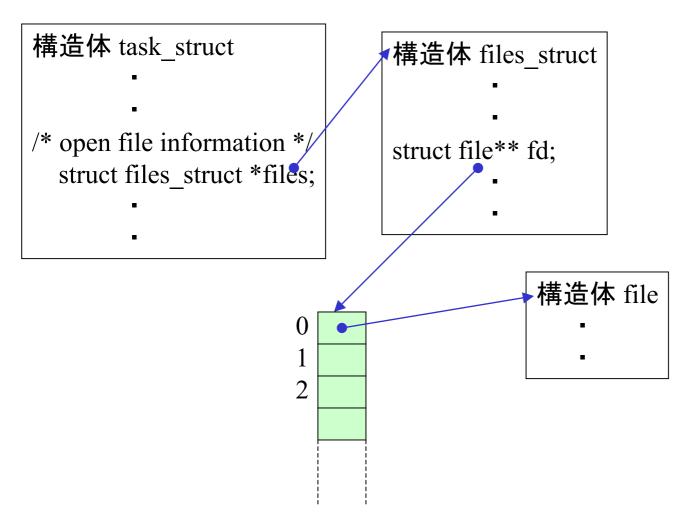
#### プロセスとファイル

- プロセスからはファイルシステムの一部(全体かもしれない)が参照可能となっている.
- 参照可能なファイルシステムの最上位をルートディレクトリと呼ぶ。
- プロセス毎に、現在、注目しているディレクトリというのは異なる。
- このようなディレクトリをCurrent Working Directory と呼ぶ.

### File Descriptor

- あるプロセスが参照しているファイル(特殊ファイルを含む)は、File Descriptorと呼ばれる数値で識別される.
- 番号, 0, 1, 2 は全プロセス共通に標準入力, 出力, エラーに対応.
- ・この数値は構造体fileの配列の添え字になっている.
  - 構造体 fileについては後述.

## 構造体の関係



# open & fopen

- 普段, fopen関数なんかで開けてるFILEなんかも, 結局は, 前述のFile Descriptorを参照している.
- openシステムコールやread,writeシステムコールはFile Descriptorを直接使ってファイルの操作を行う。

### ファイルの種類

- 通常ファイル
- ・ディレクトリ
- ・シンボリックリンク
- ブロック型デバイスファイル
- キャラクタ型デバイスファイル
- パイプ, 名前付きパイプ
- ・ソケット

この辺は本日 扱わない

### inode番号とリンク

- ファイル名というのは、割と簡単につけか えられるので、実はファイルの識別子には あまりなっていない。
- ファイルはinodeというデータ構造で内部的には管理されており、inode番号がファイルを識別するidと考えることができる。

include/linux/fs.h の338行目辺りで定義

# (ハード)リンク

- ファイルを特定ディレクトリ下に名前をつけて所属させるのがリンクである。
- 前述のように、ファイルの実体はinodeで管理されているので、
  - 一同一実体のファイルを、異なるディレクトリに 異なる名前で所属させることができる。

```
      Image: Image of the process of the
```

### シンボリックリンク

- ハードリンクには大きな制限がある.
  - 一同じファイルシステムに所属するファイルにし か適用できない。
  - 一般ユーザーはディレクトリへのリンクをはれない。
- この制限を克服するためにシンボリックリンクが作られた。
  - 存在しないファイルにさえリンクが張れる.

#### ファイルのアクセス権

- ユーザー、同ーグループ、その他で分類
- ・ モードはそれぞれ rwx
- ・他3つほどのアクセス権情報を持つ.

#### 例

```
[kaiya@linux2001 ~]% ll a.*

-rw-r--r-- 1 kaiya software 94 Jul 9 2002 a.c

-rw-r--r-- 1 kaiya software 432 Jun 29 2002 a.html

-rw-r--r-- 1 kaiya software 11226 Jun 15 02:37 a.jpg

-rwxr-xr-x 1 kaiya software 12704 Jul 29 2002 a.out*

[kaiya@linux2001 ~]% []
```

### 個々のファイル内のアクセス法

- 順次的(sequential)アクセス
  - Cのリスト構造やビデオテープのように前から順番にアクセスする方式。
  - UNIX系の通常ファイルは通常このアクセス形式をとる。
- ・ランダムアクセス
  - Cの配列のように順番に関係なく任意の位置 のデータにアクセスできる方式。

#### 通常ファイルのためのシステムコール

- open アクセスするためにファイルを開ける関数.
- close 逆に閉じる関数
- read 読む関数
- write 書く関数
- Iseek 現在の読み書き位置を移動させる関数. (ビデオの早送り巻き戻しみたいなもの)
- rename ファイル名の付け替え
- unlink ディレクトリからのエントリ削除.

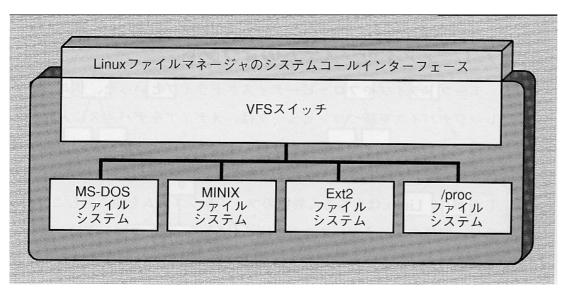
これらはシステムコールなので、実際には、カーネルへの作業要求依頼である.

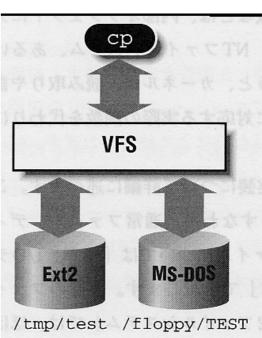
**17** 

### 仮想ファイルシステム VFS

- Linuxでは物理的・論理的に種類の異なるファイルシステムを統一的に扱うために、
- 仮想ファイルシステム(スイッチ)VFSというのを仲介役としておいて全てのファイルを扱うようにしている。
- これによって、ファイルシステムの種類や 所在をユーザーは意識せずにファイルに アクセスできる。
  - システムコールの種類を変える必要が無い.

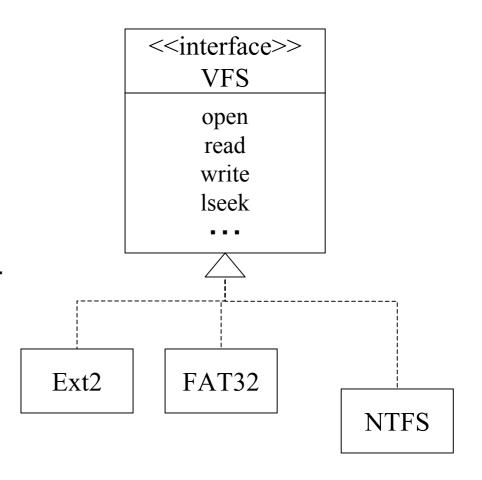
### VFSの概念図





#### VFS/Linterface

- VFSは(UML等 でいうところ の)Interfaceに あたる.
  - Linuxはオブジェクト指向言語で書いてあるのではないので概念上の話.



#### VFSが扱うファイルシステムの種類

- ディスクベースのファイルシステム
  - Ext2, FAT, NTFS, ISO9660 等のマシンに直接接続された装置内のファイルシステム群
- ネットワークファイルシステム
  - NFSやSMBなどのネットワークファイルシステム
- 特殊ファイルシステム
  - /proc/の下のファイルや /dev/pts等

#### 特殊ファイルの例 端末がファイルに?!

送り込まれた文字が端末に出 てくる。 [kaiya@linux2001 ~]% tty /dev/ttup0 [kaiya@linux2001 ~]% this is a pen. ● 潔 tcsh [kaiya@linux2001 ~]% tty /dev/ttyp1 [kaiya@linux2001 ~]% cat > /dev/ţtyp0 ≬this is a pen. [kaiya@linux2001 ~]% [] ファイルと同様に端末に対して リダイレクションすると.

### ファイルアクセスの実際

- 種類の違うファイルシステムへのアクセスには当然, 異なる関数(機能)が必要である.
  - CDROMとHDとではデータの読み方は違うだろう.
- しかし、OS操作やアプリ上において、一般にはファイルシステムの実体の違いにより関数を使い分けすることは無い。
  - CDROMだろうが、HDだろうが、openで開けてreadで読む.
- ユーザーはVFSを経由してファイルアクセスするので、こ の点を考慮しなくてよい。
- 当然, VFSは一般的なファイルアクセス関数群(システムコール)とシステム固有のアクセス関数との対応を知っている.

#### システム固有の関数までの道の例

- 1. ユーザーは**read()**等のシステムコールを呼 ぶ.
- 2. カーネルはsys\_read()を呼ぶ.
- 3. この関数は、カーネル内に保持されている開き 済なファイルを示す構造体 fileのインスタン スを読む.
- 4. fileのメンバーに構造体f\_opというのがあり、 これに関数へのポインタが列挙されている.
- 5. f\_op内のこのポインタ群がシステム固有の関数である.

file ->f\_op->read(...) という間接的な呼び方

### 関数へのポインタ

- ・ 普通のCの構文の一種.
- 返り値と引数の型が一致している関数を、 変数(コレが関数へのポインタ)に代入して、 関数呼び出しを動的に変更することができる。
- ・ 無論, 関数引数にも使える.

#### 例

```
return a+b;
int sub(int a, int b){
 return a-b;
main(){
int (*fp)(int, int); // ここで関数へのポインタを宣言.
 fp=add; // 関数名を代入できる
 printf("%d¥n", (*fp)(10, 4));
 fp=sub; // 関数名を代入できる
```

printf("%d¥n", (\*fp)(10, 4));

int add(int a, int b){

命令文は全く同じだが, fpの指している関数の実 体が異なるため,計算結 果も異なる.

## システム固有の関数までの道

- 1. ユーザーは**read()**等のシステムコールを呼 ぶ.
- 2. カーネルは**sys\_read()**を呼ぶ.
- 3. この関数は、カーネル内に保持されている開き 済なファイルを示す構造体 fileのインスタン スを読む.
- 4. fileのメンバーに構造体f\_opというのがあり、 これに関数へのポインタが列挙されている.
- 5. f\_op内のこのポインタ群がシステム固有の関数である.

file ->f\_op->read(....) という間接的な呼び方

### ファイル操作切り替えの戦略

本質的な戦略 の説明 実際はモちっと 複雑

```
ssize_t read_ext2(struct file *, char *, size_t, loff_t *) {
    // ext2ファイルシステムを読む処理.
}
ssize_t read_fat32(struct file *, char *, size_t, loff_t *) {
    // FAT32ファイルを読む処理.
}
```

```
main() {
ssize_t (*read)(struct file *, char *, size_t, loff_t *);

if(ファイルがext2上にある)
read=read_ext2;
else if(ファイルがFAT32上にある)
read=read_fat32;

(*read)(&file, buf, size); // 実体に関係なくファイルを読む
}
```

### 構造体 file

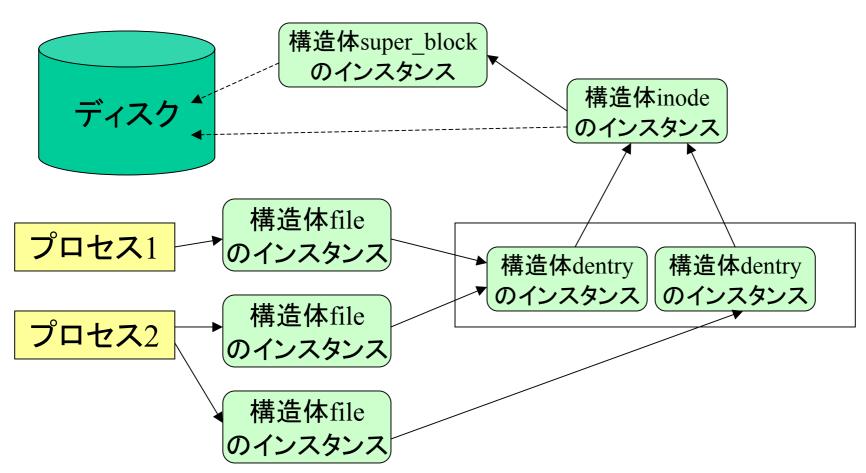
include/linux/fs.h struct file { の423行目くらい struct file \*f next, \*\*f pprev; struct dentry \*f dentry; struct file\_operations \*f\_op; // コイツがInterfaceにあたる f mode; mode t loff t f pos; unsigned int f count, f flags; unsigned long f reada, f ramax, f\_raend, f\_ralen, f\_rawin; struct fown struct f owner; unsigned int f uid, f gid; int f error; unsigned long f version; /\* needed for tty driver, and maybe others \*/ void \*private data;

# 構造体 file\_operations

```
struct file operations {
                                                                include/linux/fs.h
  loff t (*llseek) (struct file *, loff t, int);
                                                                の810行目くらい
  ssize t (*read) (struct file *, char *, size t, loff t *);
  ssize t (*write) (struct file *, const char *, size t, loff t *);
  int (*readdir) (struct file *, void *, filldir t);
  unsigned int (*poll) (struct file *, struct poll table struct *);
  int (*ioctl) (struct inode *, struct file *, unsigned int, unsigned long);
  int (*mmap) (struct file *, struct vm_area_struct *);
  int (*open) (struct inode *, struct file *);
  int (*flush) (struct file *);
  int (*release) (struct inode *, struct file *);
                                                          メンバが関数への
  int (*fsync) (struct file *, struct dentry *);
                                                          ポインタとなってお
  int (*fasync) (int, struct file *, int);
                                                          り、ポインタの指す
  int (*check media change) (kdev t dev);
                                                          実体がFSの実体毎
  int (*revalidate) (kdev t dev);
  int (*lock) (struct file *, int, struct file lock *);
                                                          に異なる.
};
```

全てのファイルシステムに全ての関数が定義されているわけではない.

## プロセスとVFS間の関係



個々のプロセスがファイルシステムにたどり着くまで、いくつかの構造体のインスタンスを経由します.

31

### 各構造体の意味の説明

- file オープンされているファイルの情報を 保持する. プロセスがファイルにアクセスしている間しかインスタンスは存在しない.
- dentry ディレクトリのエントリを現すインスタンス.
- inode 個々のファイルについての一般的情報.
- super\_block マウントされたファイルシステムに関する情報を保持.

## 構造体 dentry

```
struct dentry {
  int d count;
  unsigned int d flags;
  struct inode * d inode; /* Where the name belongs to - NULL is negative$
  struct dentry * d parent; /* parent directory */
  struct dentry * d mounts; /* mount information */
  struct dentry * d covers;
  struct list head d hash; /* lookup hash list */
  struct list head d lru; /* d count = 0 LRU list */
                                                       include/linux/dcash.h
  struct list head d child; /* child of parent list */
                                                       の58行目くらい
  struct list head d subdirs; /* our children */
  struct list_head d_alias; /* inode alias list */
  struct qstr d name;
  unsigned long d time; /* used by d revalidate */
  struct dentry operations *d op;
  struct super block * d sb; /* The root of the dentry tree */
  unsigned long d reftime; /* last time referenced */
  void * d fsdata; /* fs-specific data */
  unsigned char d iname[DNAME INLINE LEN]; /* small names */
};
```

### 構造体 inode

```
struct inode {
  struct list head i hash;
  struct list head i list;
  struct list head i dentry;
  unsigned long i ino;
  unsigned int i count;
 kdev t i dev;
 umode_t i mode;
 nlink_t i_nlink;
 uid_t i_uid;
  gid_t i_gid;
 kdev_t i_rdev;
 off t i size;
 time_t i_atime;
 time t i mtime;
  time t i ctime;
```

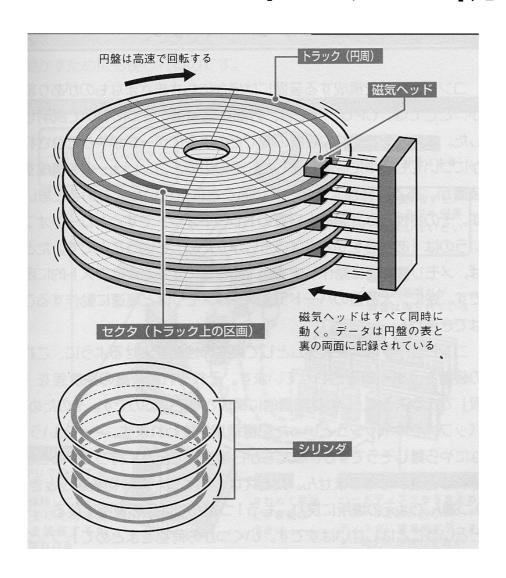
```
unsigned long i blksize;
  unsigned long i_blocks;
  unsigned long i version;
  unsigned long i nrpages;
  struct semaphore i sem;
  struct semaphore i atomic write;
  struct inode operations *i op;
  struct super block *i sb;
  struct wait queue *i wait;
  struct file lock *i flock;
  struct vm_area_struct *i mmap;
  struct page *i pages;
  struct dquot *i dquot[MAXQUOTAS];
// 長いので以下省略
```

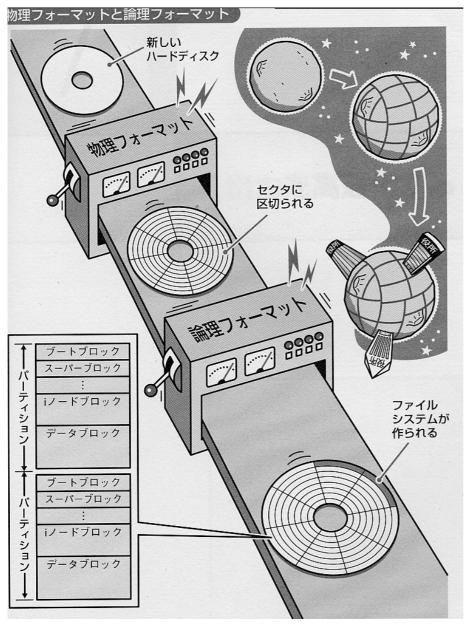
# 構造体 super\_block

```
struct super block {
  struct list head s list; /* Keep this first */
  kdev t
              s dev;
  unsigned long s blocksize;
  unsigned char s blocksize bits;
  unsigned char
                 s lock;
  unsigned char
                 s rd only;
  unsigned char
                  s dirt;
  struct file system type *s type;
  struct super_operations *s_op;
  struct dquot_operations *dq_op;
  unsigned long
                  s flags;
  unsigned long s magic;
  unsigned long
                   s time;
  struct dentry
                  *s root:
  struct wait queue *s wait;
  struct inode
                 *s ibasket:
  short int
             s ibasket count;
             s ibasket_max;
  short int
  struct list head s dirty; /* dirty inodes */
 // 以下長いので略
```

include/linux/fs.h の528行目くらい

# ディスクの物理構造





#### 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。

## 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	Linux
4つのパー 📗	/dev/hda2	256	321	530145	82	Linux スワップ
ティション	/dev/hda3	322	576	2048287+	83	Linux
, , , , , ,	/dev/hda4	577	2434	14924385	83	システム Linux Linux スワップ Linux Linux

