

# Journaling Filesystem

**VA Linux Systems Japan**

岡島順治郎

**jro@valinux.co.jp**

Copyright(c) 2002, VA Linux Systems Japan K.K. All rights reserved.

Linux World VA Linux Kernel Forum <2002/05/29>



# ファイル生成時間を計る time(bash)

```
#!/bin/bash -
n=${1:-100}
while test $n -gt 0; do s="$s $n"; let n-=1; done
time touch $s
rm $s
```

	<b>ext2fs</b>	<b>Another ext2fs</b>
<b>real</b>	<b>0m0.055s</b>	<b>0m2.184s</b>
<b>user</b>	<b>0m0.000s</b>	<b>0m0.000s</b>
<b>sys</b>	<b>0m0.010s</b>	<b>0m0.000s</b>

# ファイル生成時間を計る iostat(1) - 1

```
#!/bin/bash -
n=${1:-100}
while test $n -gt 0; do s="$s $n"; let n-=1; done
set -m
iostat -d 1 &
touch $s
sleep 1
kill %iostat
rm $s
```

# ファイル生成時間を計る iostat(1) - 2

<b>ext2fs</b>	<b>tps</b>	<b>Bread/s</b>	<b>Bwrite/s</b>	<b>Bread</b>	<b>Bwrtn</b>
	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	<b>0</b>	<b>0</b>
<b>Another ext2fs</b>	<b>tps</b>	<b>Bread/s</b>	<b>Bwrite/s</b>	<b>Bread</b>	<b>Bwrtn</b>
	<b>120.00</b>	<b>0.00</b>	<b>960.00</b>	<b>0</b>	<b>960</b>
	<b>115.00</b>	<b>0.00</b>	<b>928.00</b>	<b>0</b>	<b>928</b>
	<b>113.00</b>	<b>0.00</b>	<b>904.00</b>	<b>0</b>	<b>904</b>
	<b>116.00</b>	<b>0.00</b>	<b>936.00</b>	<b>0</b>	<b>936</b>

# Mount option, sync と async

touch(1)の例からわかること

**sync option**

open(O\_CREAT)の最中に何かを大量に書き込み、時間がかかる。

**async option**

基本的に書き込みを行わないため、時間がかからない。

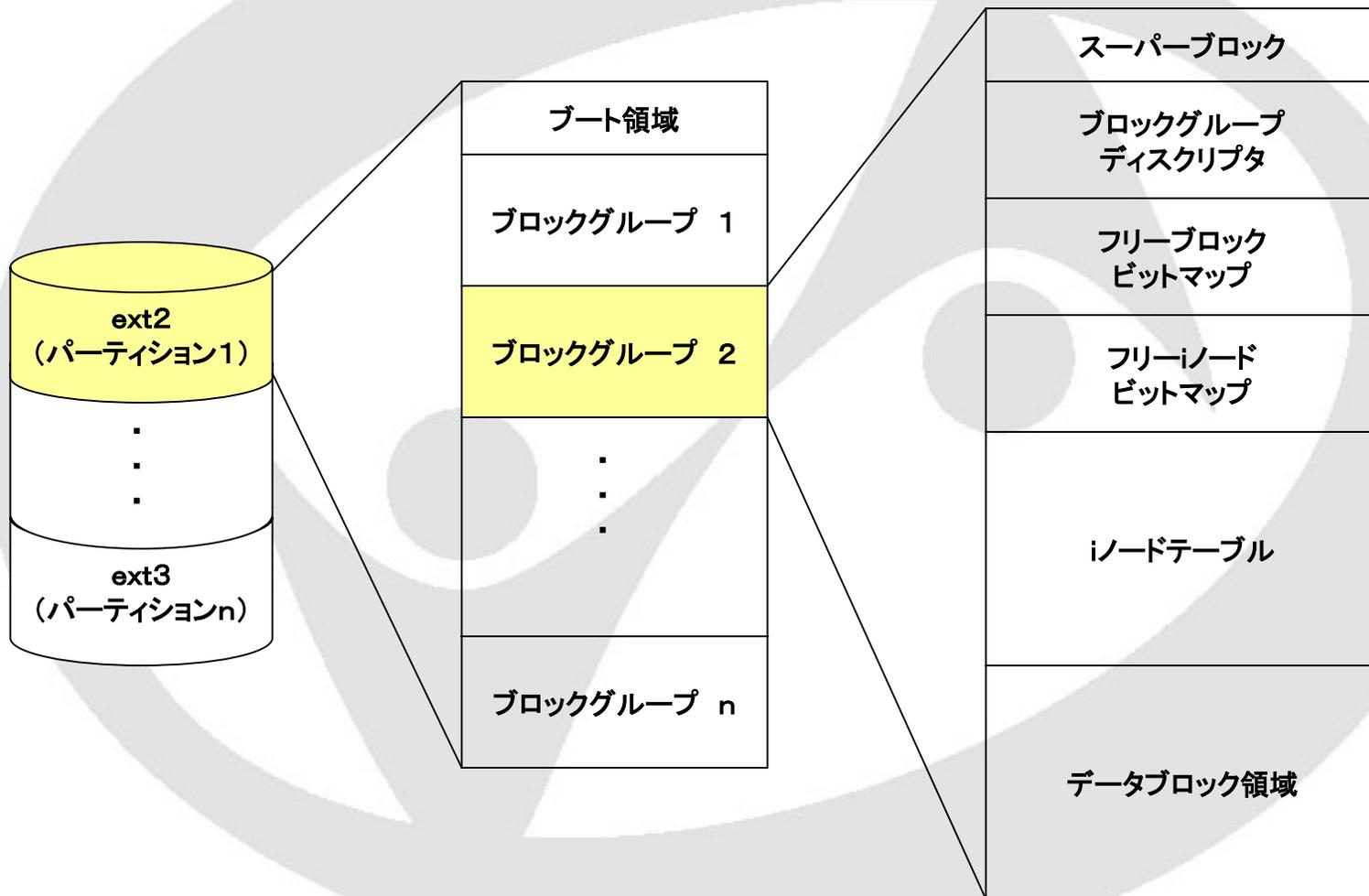
<shameless plug>

<a href="http://www03.u-page.so-net.ne.jp/da2/h-takaha/">

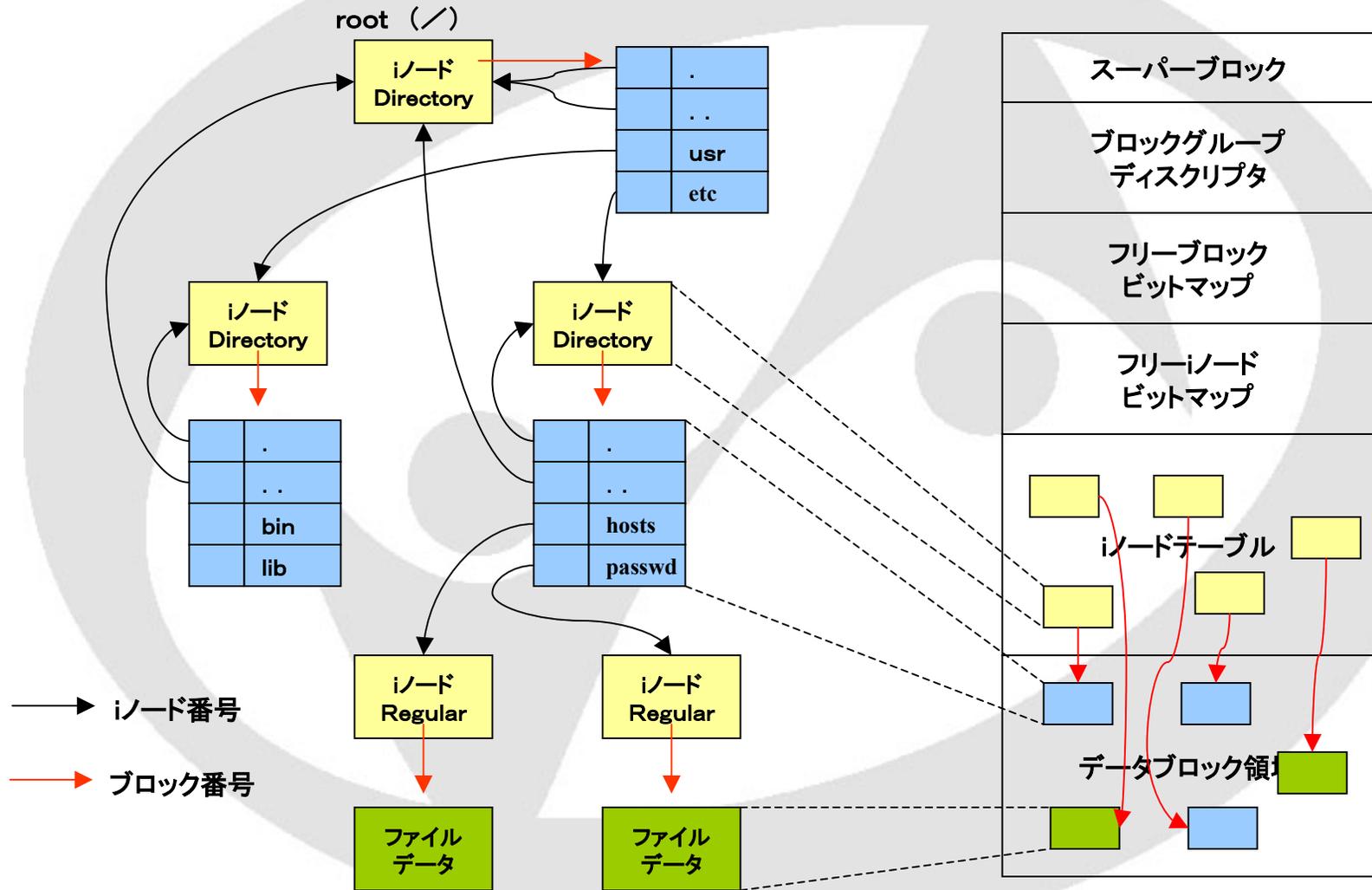
**open(2)の  
O\_SYNC**

open(2)後のwrite(2)に作用する。touch(1)によるファイル作成ではファイルの中身(ファイルデータ)は書き込んでいない。

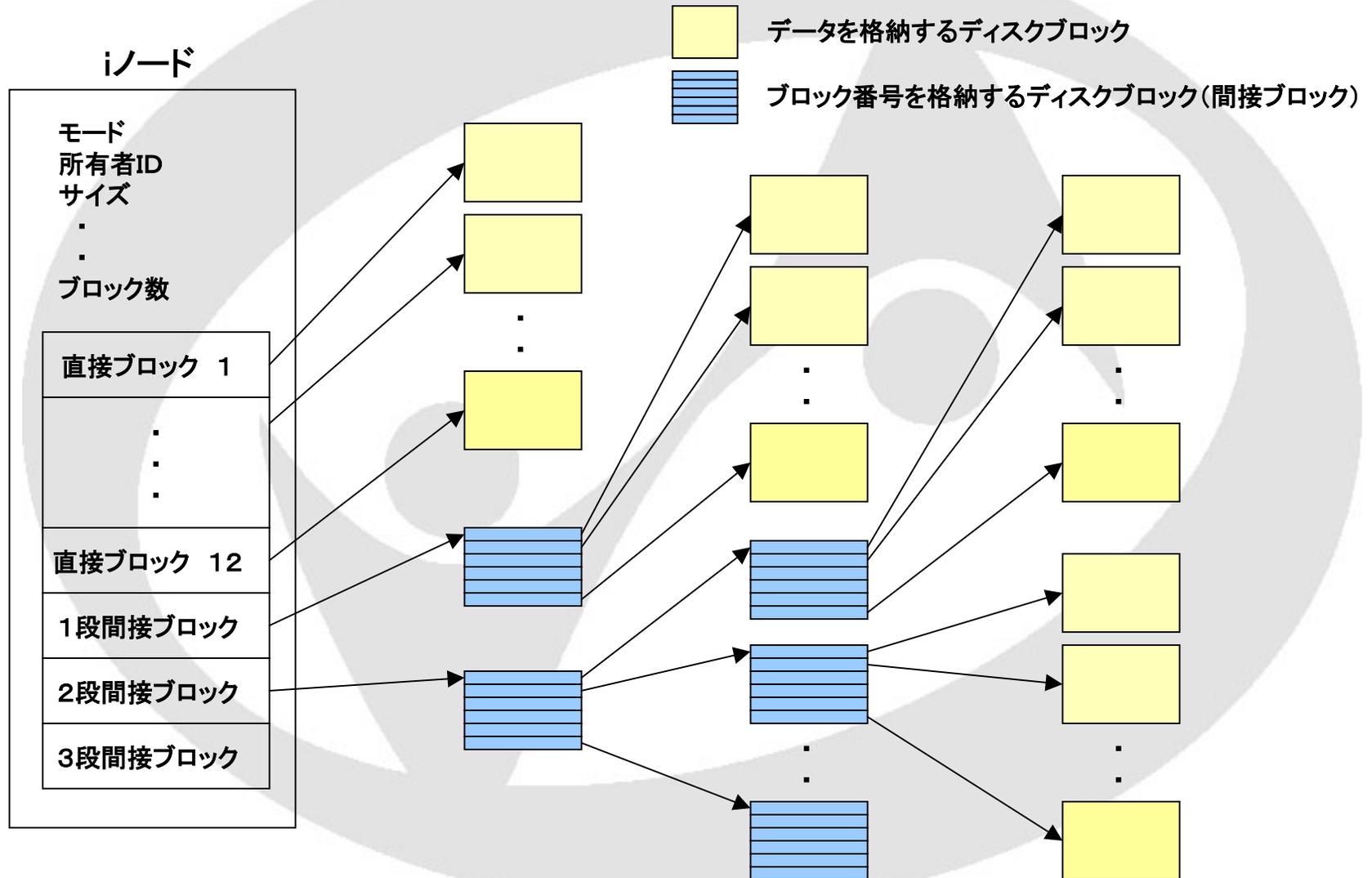
# ext2:レイアウト



# ext2: ディレクトリ階層



# ext2: ファイルデータ構造



Copyright(c) 2002, VA Linux Systems Japan K.K. All rights reserved.

# メタデータ

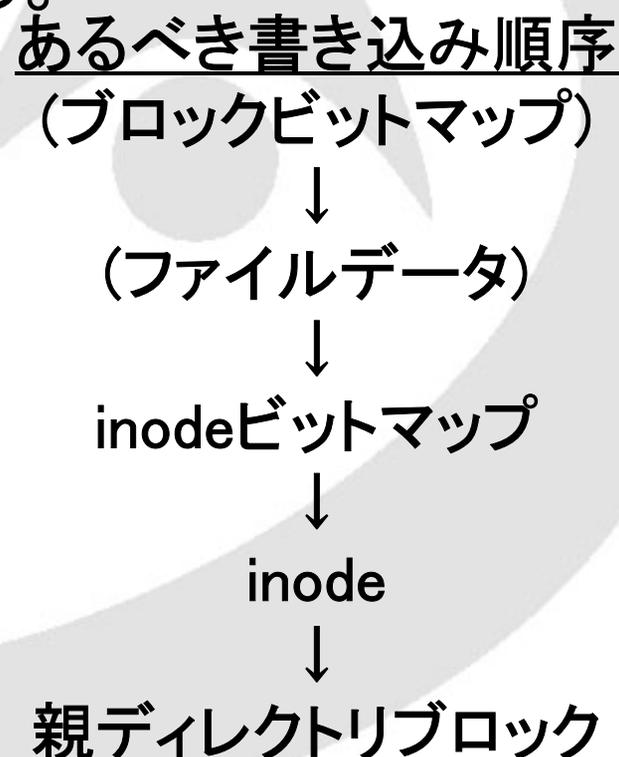
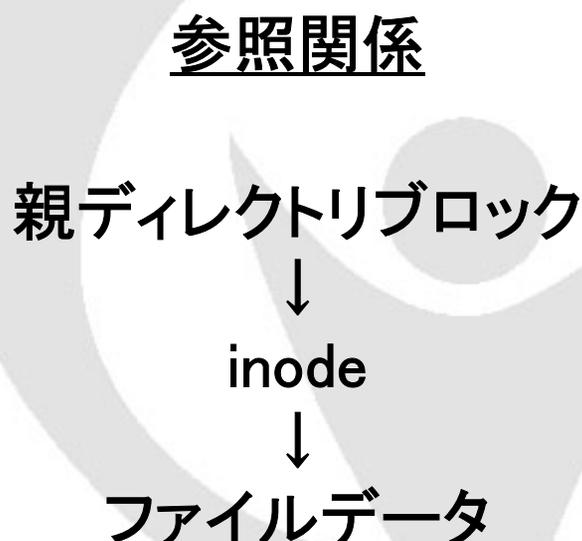
ファイルシステム上でファイルの中身(ファイルデータ)以外の情報。

touch(1)の例から取り上げるものは、以下の三つ。

親ディレクトリ ブロック	ファイル名とinode番号の対が書き込まれている。
inode	ファイルに関するさまざまな情報。ファイルデータがある場合は、そのブロック番号など。touch(1)の例では、ファイルデータは無い。
inodeビットマップ	ファイルシステム上のinode達が使用中・未使用かを表す。

# メタデータの書き込み順序

書き込み順序は参照関係に従うのが良い。  
書き込まれていないものを参照しない状態を常に保てば、障害耐性の向上につながる。



# メタデータ更新の実例 creat("/d/f") -- 1

Output of dumpe2fs	Block group 0	Block group 1
Superblock	1(block number)	8193
G. descriptor	2	8194
Block bitmap	3	8195
inode bitmap	4	8196
inode table	5 - 260	8197 - 8452

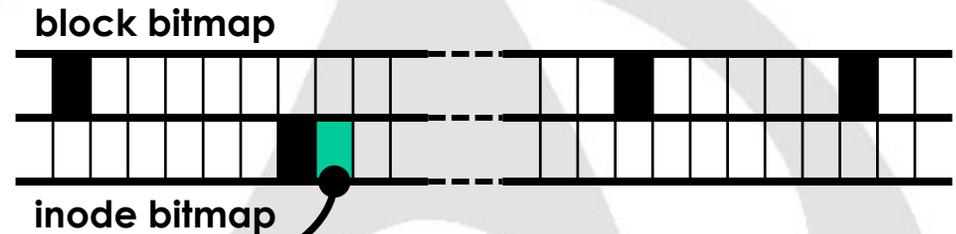
Output of debugfs	/	/d	/d/f
inode	2(inode number)	2049	2050
Type	Directory	Directory	Regular
Blocks	261(block number)	8453	8454

# メタデータの参照関係 creat("/d/f")--2

output of "debugfs"

ブロック <u>261</u>	
rootディレクトリブロック	
2	.
2	..
<u>2049</u>	d
....	....

ブロック <u>8453</u>	
/dディレクトリブロック	
2049	.
2	..
<u>2050</u>	f



ブロック <u>8197</u>	
inode table	
inode <u>2049</u>	/d inode
データブロック <u>8453</u>	
inode <u>2050</u>	/d/f inode
データブロック <u>0</u>	

# 非同期書き出し

Systemcall発行時にdiskへの書き出しは行われず、  
unmount時、kernel threadのbdflushやkupdatedにより  
書き込まれる。

**bdflush**

dirty buffer が一定以上た  
まったら動きだし、書き出し  
を行い、dirty buffer 数を  
一定数まで減らす

**kupdated**

一定間隔をおいて起動し、  
一定時間を経過したdirty  
buffer を書き出す

# 整合性とfsck

---

## ファイルシステムの整合性

- `kupdated`, `bdflush`により、バッファが書き出される前
- バッファが書き出された後
- バッファが書き出されている最中

## 予期せぬシステムダウン後のfsck

- 修復の方針
- 所要時間

# ext2 + jbd = ext3(1)

---

## ファイルシステムの整合性

- 従来のkupdated, bdflushなどの動作を変えずに、メタデータなどの更新情報(ログ)を独立した領域に、五秒間隔で書き込む。

# ext2 + jbd = ext3(2)

---

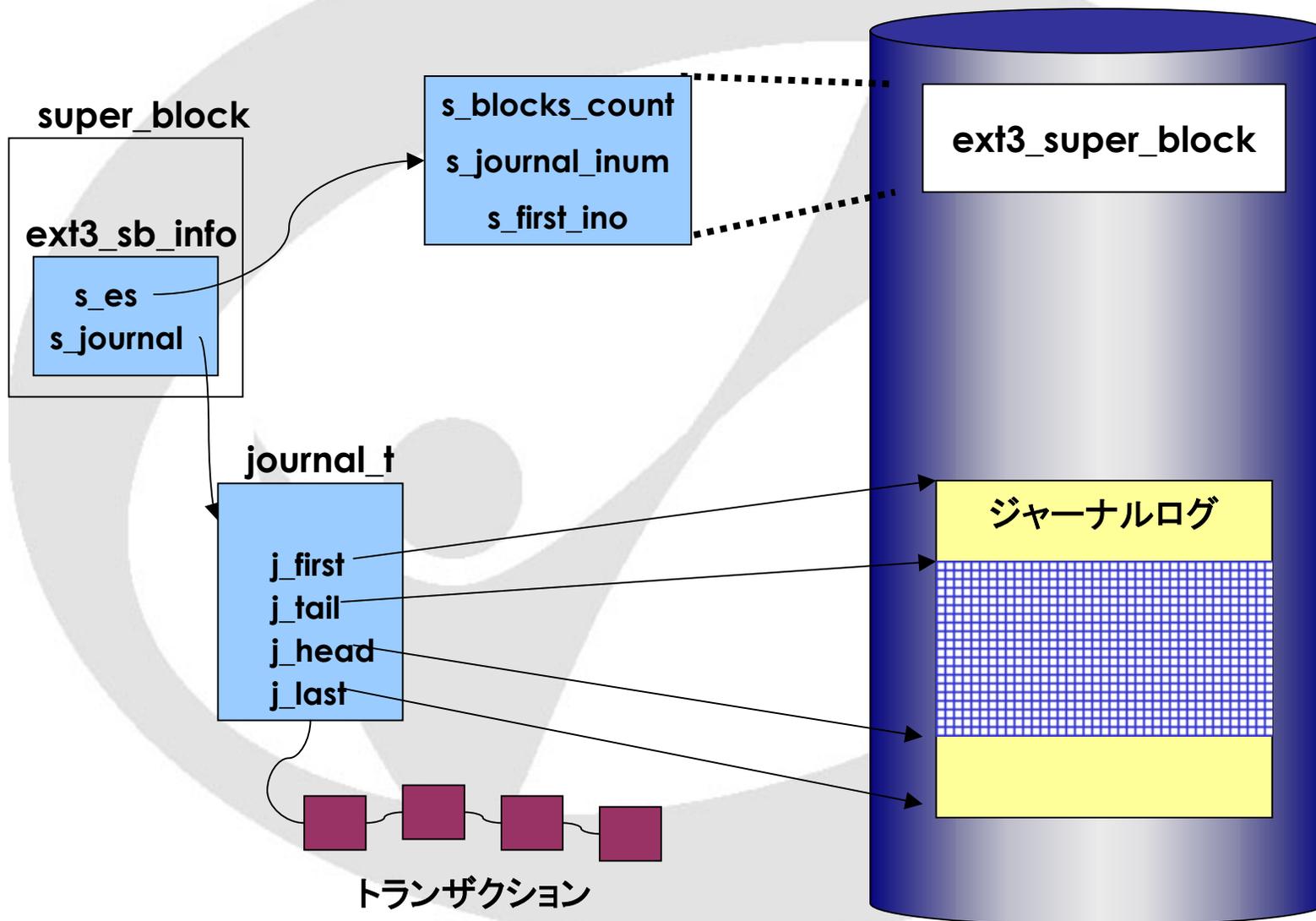
予期せぬシステムダウン後のfsck

- ジャーナルログからまるごと復旧できる。
- 対象範囲がジャーナルログだけなので、所要時間が非常に短い。

特徴

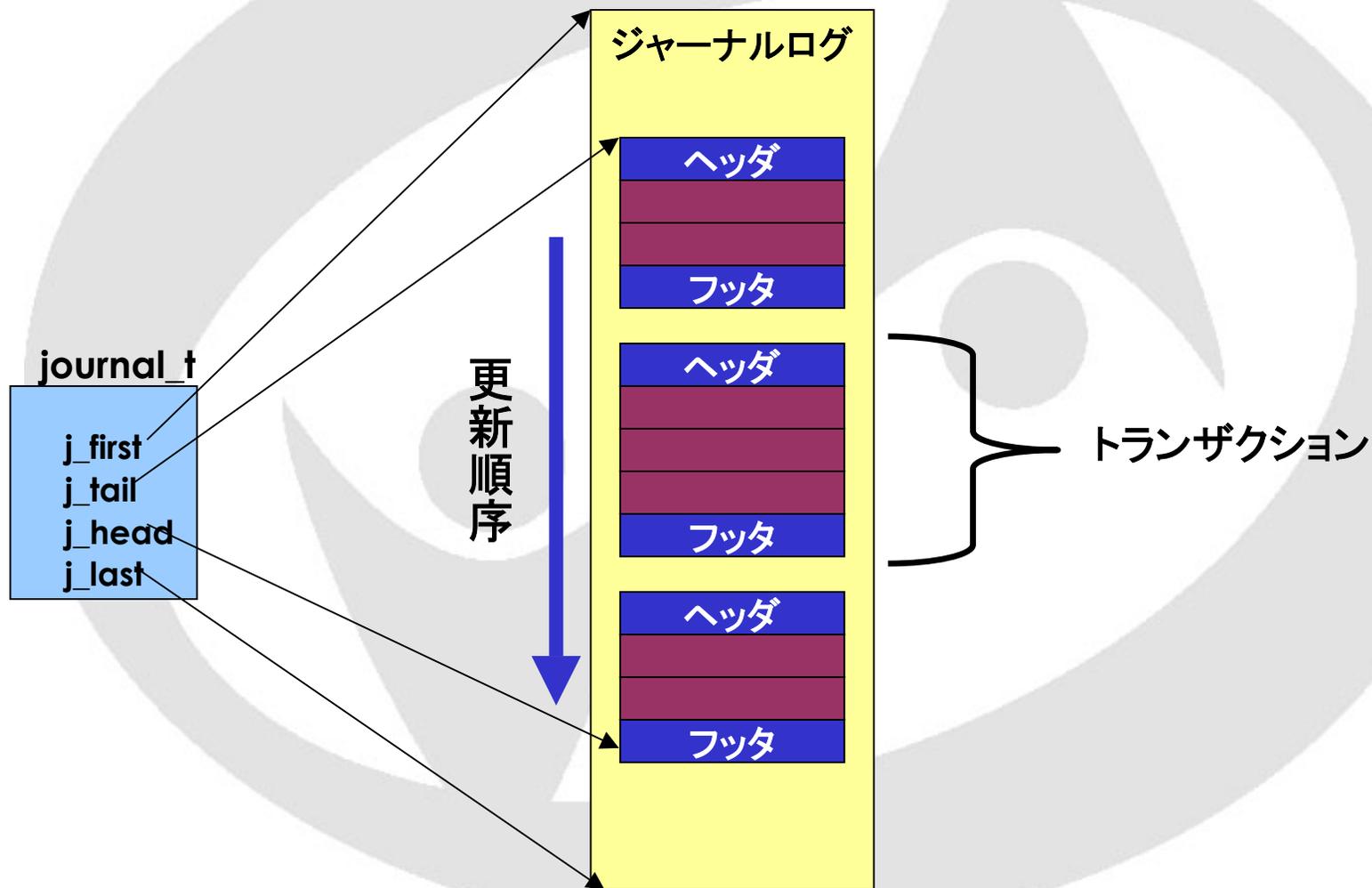
- ディスクレイアウトがext2fsと互換性がある。ジャーナル領域を予約inodeにつなぎ、ファイルデータとしている。
- jbdはファイルシステムに依存しておらず、汎用性と独立性が高い。

# EXT3ジャーナルデータ構造

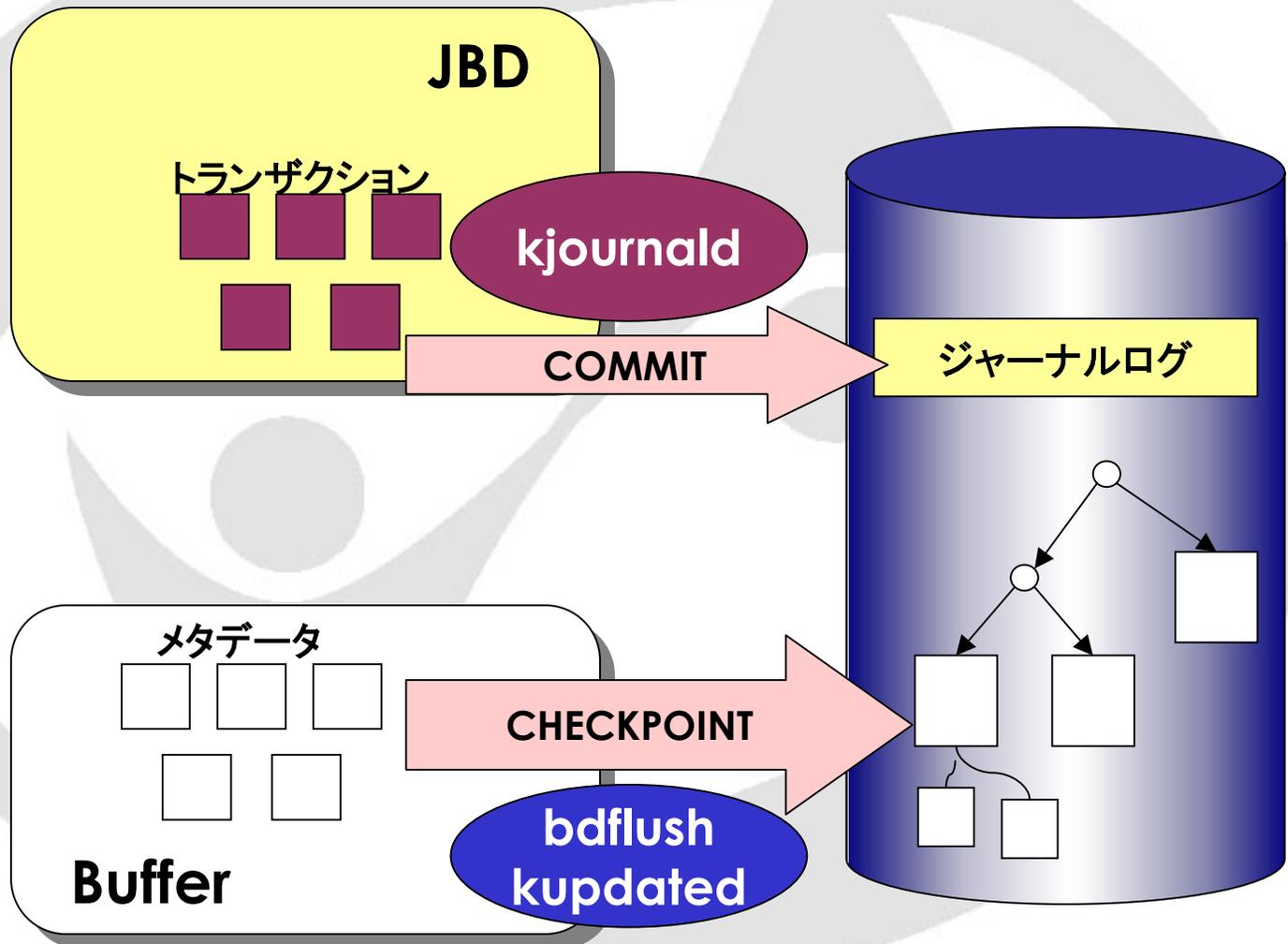


Copyright(c) 2002, VA Linux Systems Japan K.K. All rights reserved.

# ジャーナルログ形式



# ジャーナルのCOMMITとCHECKPOINT



# ファイル生成時間を計る time(bash)

```
#!/bin/bash -
n=${1:-100}
while test $n -gt 0; do s="$s $n"; let n-=1; done
time touch $s
rm $s
```

	ext3fs		ext2fs	
	async	sync	async	sync
<b>real</b>	<b>0m0.364s</b>	<b>0m9.364s</b>	<b>0m0.055s</b>	<b>0m2.184s</b>
<b>user</b>	<b>0m0.010s</b>	<b>0m0.010s</b>	<b>0m0.000s</b>	<b>0m0.000s</b>
<b>sys</b>	<b>0m0.300s</b>	<b>0m0.300s</b>	<b>0m0.010s</b>	<b>0m0.000s</b>

※非常に単純な測定！鵜呑みにするな。