



**TECHNISCHE  
UNIVERSITÄT  
DRESDEN**

**Department of Computer Science** Institute of System Architecture, Operating Systems Group

# **EIO: ERROR CHECKING IS OCCASIONALLY CORRECT**

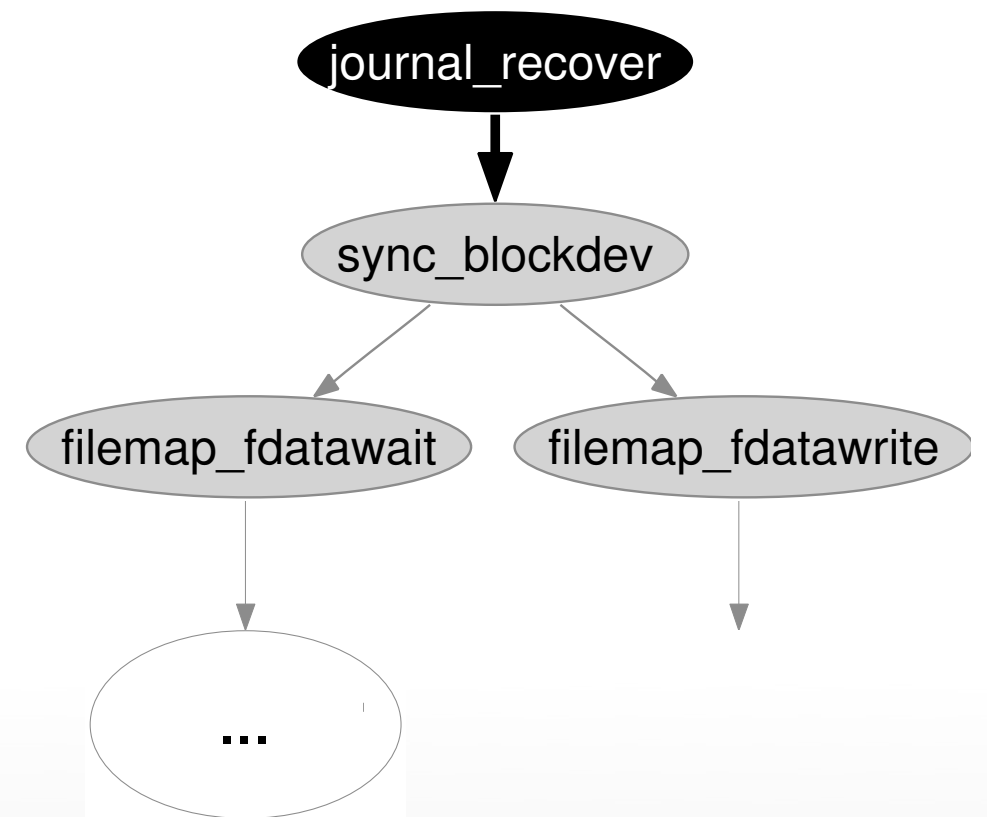
**HARYADI S. GUNAWI, CINDY RUBIO-GONZÁLEZ, ANDREA C.  
ARPACI-DUSSEAU, REMZI H. ARPACI-DUSSEAU, BEN LIBLIT**

**CARSTEN WEINHOLD**

- File and storage systems must be robust
- **Previous research:** *„file systems are [...] unreliable when the underlying disk system does not behave as expected“*
- **Requirement:** comprehensive recovery policies need correct error reporting
- **Reality:** error propagation often incorrect
- Paper presents analysis error propagation in Linux code

- Error Detection and Propagation (EDP):
  - Static analysis of dataflow (error codes)
  - Uses source-to-source transformation
  - Tracks error propagation through call stacks
- Used to analyze Linux 2.6.15.4 source:
  - VFS, memory management
  - All file systems (ext3, XFS, NFS, VFAT, ...)
  - SCSI, IDE, soft RAID storage subsystems

- Basic abstraction: channels
  - Set of function calls
  - **Generation endpoint:**  
error first exposed
  - **Termination endpoint:**  
end of error propagation
  - **Propagating functions** in  
between



```
struct file_ops {  
    int (*read) ();  
    int (*write) ();  
};
```

```
switch (...) {  
    case ext2:  ext2_read();  break;  
    case ext3:  ext3_read();  break;  
    case ntfs:  ntfs_read();  break;  
    ...  
}
```

```
struct file_ops ext2_f_ops {  
    .read  = ext2_read;  
    .write = ext2_write;  
};  
struct file_ops ext3_f_ops {  
    .read  = ext3_read;  
    .write = ext3_write;  
};
```

$\exists$  *if (expr) { ... }, where  
errorCodeVariable  $\subseteq$  expr*

- Error-complete channels:

```

1 void goodTerminationEndpoint() {
2     int err = generationEndpoint();
3     if (err)
4         ...
5 }
6 int generationEndpoint() {
7     return -EIO;
8 }

```

- Error-broken channels:

```

1 void badTerminationEndpoint() {
2     int err = generationEndpoint();
3     return;
4 }

```

**Unchecked  
error**

```

1 // hfs/bfind.c
2 int find_init(find_data *fd) {
3     fd->search_key = kmalloc(...);
4     if (!fd->search_key)
5         return -ENOMEM;
6     ...
7 }
8 // hfs/ilode.c
9 int file_lookup() {
10    find_init(fd); /* NOT-SAVED E.C */
11    fd->search_key->cat = ...; /* BAD!! */
12    ...
13 }

```

**Unsaved error  
"Bad call"**

- Bad calls not always bad:
  - Multiple error returned, check only one
  - Rely on other callees to check errors

```
1 // fs/buffer.c
2 int sync_dirty_buffer (buffer_head* bh) {
3     ...
4     return ret; // RETURN ERROR CODE
5 }
6 // reiserfs/journal.c
7 int flush_commit_list() {
8     sync_dirty_buffer(bh); // UNSAVED EC
9     if (!buffer_uptodate(bh)) {
10         return -EIO;
11     }
12 }
```

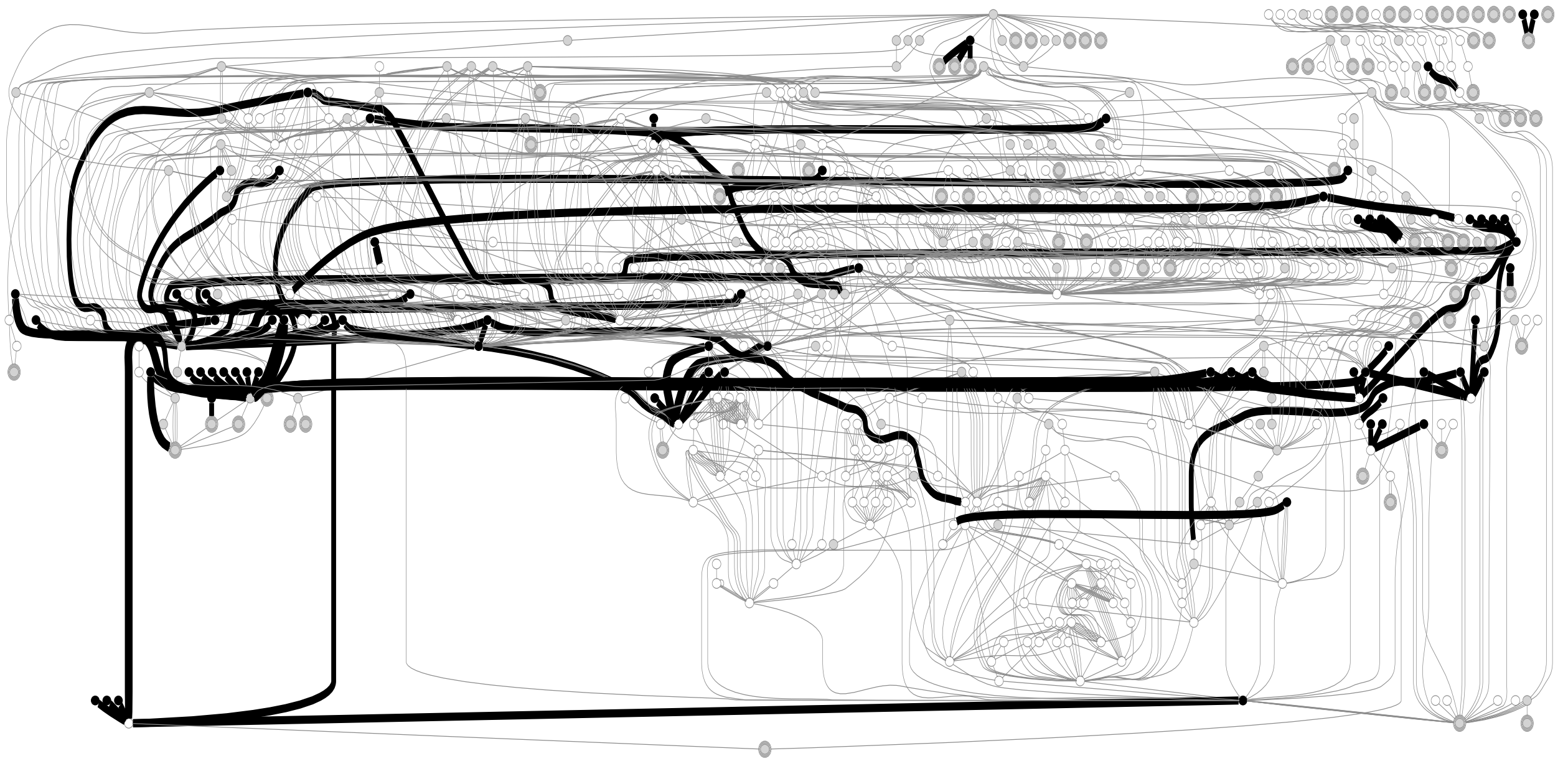




# EXAMPLE: HFS

<b>Viol#</b>	<b>Caller</b> → <b>Callee</b>	<b>Filename</b>	<b>Line#</b>
<b>A</b>	file_lookup → find_init	inode.c	493
<b>B</b>	fill_super → find_init	super.c	385
<b>C</b>	lookup → find_init	dir.c	30
<b>D</b>	brec_updt_prnt → __brec_find	brec.c	405
<b>E</b>	brec_updt_prnt → __brec_find	brec.c	345
<b>F</b>	cat_delete → free_fork	catalog.c	228
<b>G</b>	cat_delete → find_init	catalog.c	213
<b>H</b>	cat_create → find_init	catalog.c	95
<b>I</b>	file_trunc → free_exts	extent.c	507
<b>J</b>	file_trunc → free_exts	extent.c	497
<b>K</b>	file_trunc → find_init	extent.c	494
<b>L</b>	ext_write_ext → find_init	extent.c	135
<b>M</b>	ext_read_ext → find_init	extent.c	188
<b>N</b>	brec_rmv → __brec_find	brec.c	193
<b>O</b>	readdir → find_init	dir.c	68
<b>P</b>	cat_move → find_init	catalog.c	280
<b>Q</b>	brec_insert → __brec_find	brec.c	145
<b>R</b>	free_fork → free_exts	extent.c	307
<b>S</b>	free_fork → find_init	extent.c	301

**XFS** [ 105 bad / 1453 calls, 7% ]



Rank	By % Broken		By Viol/Kloc		Bad Calls	EC Calls	Size (Kloc)	Frac (%)	Viol/Kloc	
	FS	Frac.	FS	Viol/Kloc						
1	IBM JFS	24.4	ext3	7.2	SCSI (root)	<b>123</b>	628	198	19.6	0.6
2	ext3	22.1	IBM JFS	5.6	IDE (root)	<b>53</b>	223	15	23.8	3.5
3	JFFS v2	15.7	NFS Client	3.6	Block Dev (root)	<b>39</b>	195	36	20.0	1.1
4	NFS Client	12.9	VFS	2.9	Software RAID	<b>31</b>	290	32	10.7	1.0
5	CIFS	12.7	JFFS v2	2.2	SCSI (aacraid)	<b>30</b>	76	7	39.5	4.8
6	MemMgmt	11.4	CIFS	2.1	SCSI (lpfc)	<b>14</b>	30	16	46.7	0.9
7	ReiserFS	10.5	MemMgmt	2.0	Blk Dev (P-IDE)	<b>11</b>	17	8	64.7	1.5
8	VFS	8.4	ReiserFS	1.8	SCSI aic7xxx	<b>8</b>	62	37	12.9	0.2
9	NTFS	8.1	XFS	1.4	IDE (pci)	<b>5</b>	106	12	4.7	0.4
10	XFS	6.9	NFS Server	1.2						

- Only „complex“ file systems:  
10k+ SLOC, 50+ error related calls
- Ext3, JFS least robust, XFS most
- Storage: IDE has more violations than SCSI

- More than 63% of write errors ignored
- Possible explanations:
  - No higher-level error handling
  - Errors neglected intentionally

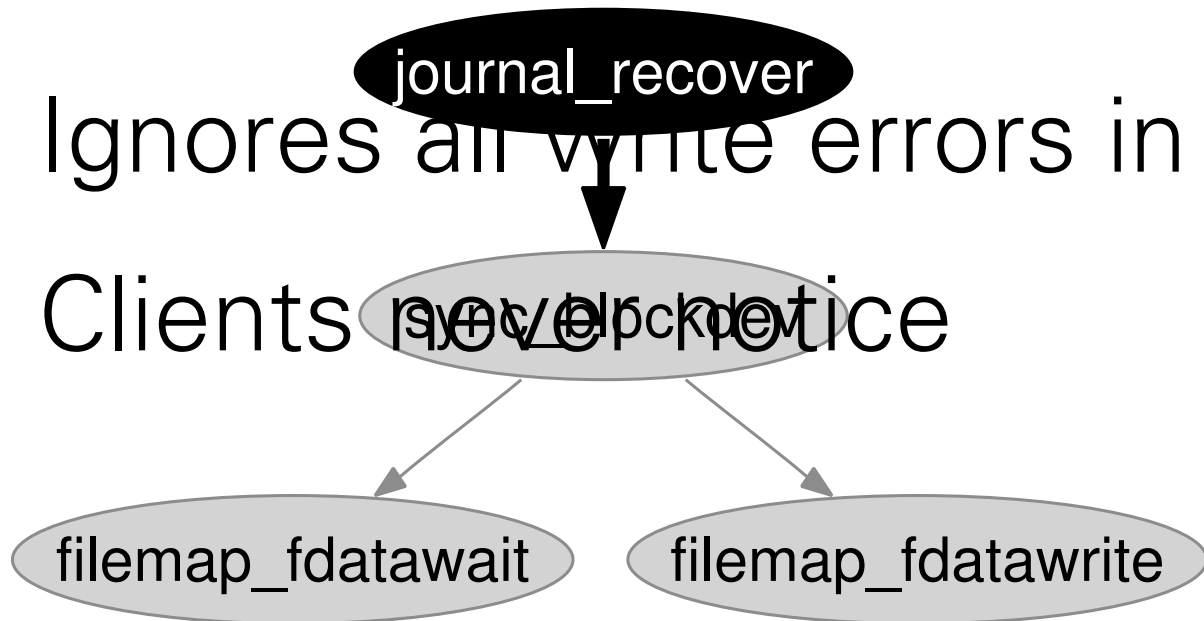
Callee Type	Bad Calls	EC Calls	Frac. (%)
Read*	26	603	<b>4.3</b>
Sync	70	236	<b>29.7</b>
Wait	27	70	<b>38.6</b>
Write	80	598	<b>13.4</b>
Sync+Wait+Write	177	904	<b>19.6</b>
Specific Callee			
filemap_fdatawait	22	29	<b>75.9</b>
filemap_fdatawrite	30	47	<b>63.8</b>
sync_blockdev	15	21	<b>71.4</b>

- **Example 1:** Journaling Block Device (JBD)

- JBD recovery code ignores all write errors
- Error code dropped in middle of channel

- **Example 2:** NFS server

- Ignores all write errors in sync writes
- Clients never notice

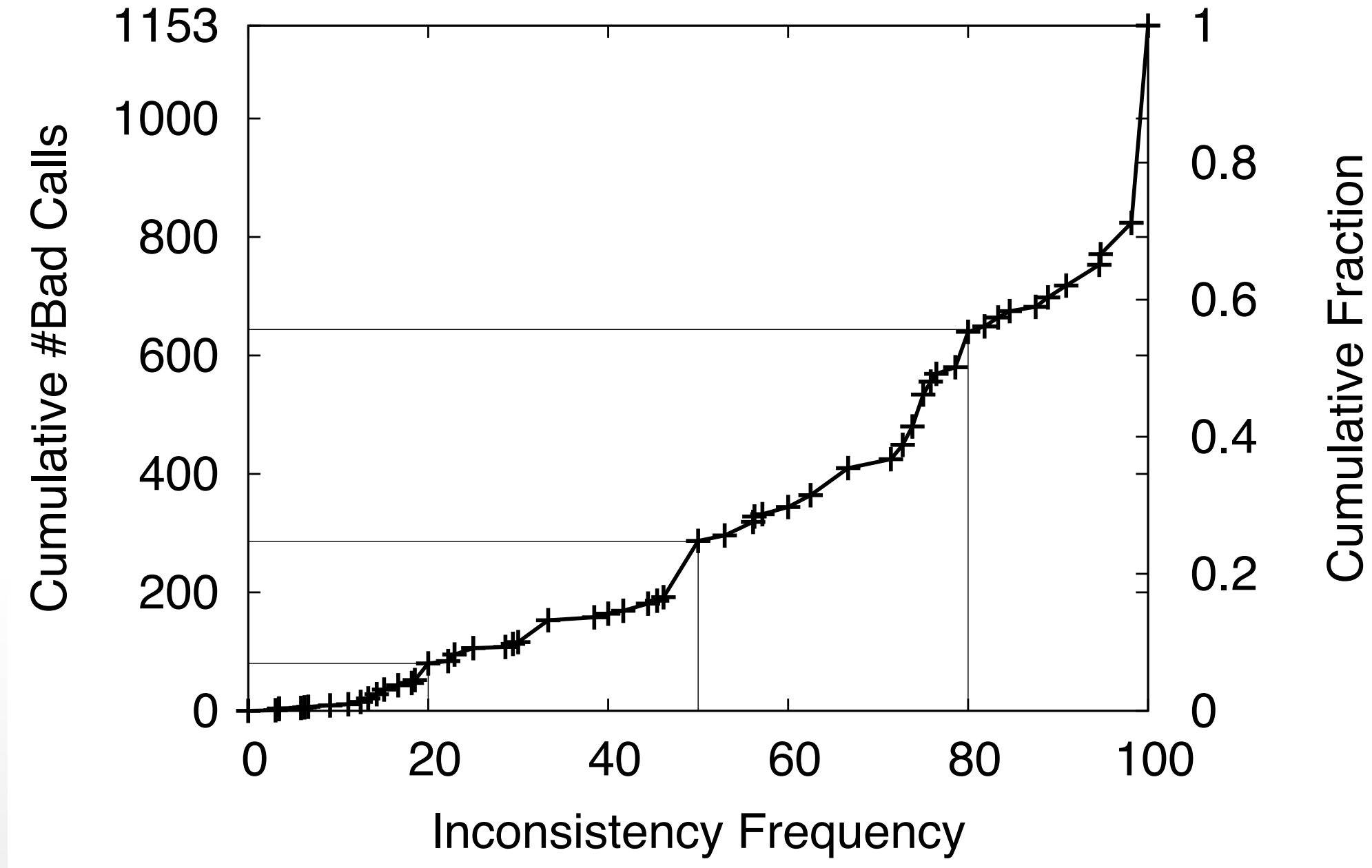


```

journal_recover()
/* BROKEN CHANNEL */
sync_blockdev();

sync_blockdev()
ret = fm_fdatawrite();
err = fm_fdatawait();
if(!ret) ret = err;
/* PROPAGATE EIO */
return ret;
  
```

CDF of Inconsistency Frequency vs. #Bad Calls



- Where are error codes dropped?
- No clear pattern:
  - File systems:
    - 10% direct, 14% later
  - Storage drivers:
    - 20% direct, 15% later

- Call distance?

	Bad Calls	EC Calls	Frac. (%)
<i>File Systems</i>			
Inter-module	307	1944	<b>15.8</b>
Inter-file	367	2786	<b>13.2</b>
Intra-file	159	2548	<b>6.2</b>
<i>Storage Drivers</i>			
Inter-module	48	199	<b>24.1</b>
Inter-file	92	495	<b>18.6</b>
Intra-file	180	1050	<b>17.1</b>



- Erros are not propagated correctly:  
Result: **1153** calls drop error (that's **13%**)
- Complex file systems are more likely to propagate errors incorrectly
- Popular file systems not the most robust
- Write errors consistently ignored:
  - May cause silent failure
  - Often no easy way to handle



- EDP catches only simple bugs, but reports many violations in all Linux file systems.
- Are the violations really that bad?
- Is OK to ignore write errors after all?
- Is ignoring write errors the *disease* or in fact a *symptom* of higher-level problems?
- Half the code is for error checking, is C the right language for that?