

PerSeVerE: Persistency Semantics for Verification Under Ext4

*Michalis Kokologiannakis*¹ Ilya Kaysin² Azalea Raad³ Viktor Vafeiadis¹

January, 2021

¹MPI-SWS

²National Research University, Higher School of Economics, Russia/JetBrains Research

³Imperial College London

Save buffer as “foo.txt”

```
d_f = open (“foo.txt”, O_WRONLY|O_CREAT|O_TRUNC);  
write (d_f, buffer);  
close (d_f);  
printf (“File saved\n”);
```

Does **replace-via-truncate** successfully write foo.txt?

- Consistency: ✓
- Persistency: ???

What Does `man` Say?

"If `auto_da_alloc` is enabled, `ext4` will detect the [...] `replace-via-truncate` [...] and [...] the data blocks of the new file are forced to disk [on close] [...]."

`man 5 ext4`

`replace-via-truncate` is used by e.g., `nano` for writing files and backups

Upon a crash, both the edited file and its backup might be empty!

Formal Semantics

- POSIX written in prose
- ext4 not fully POSIX-compliant
- ext4 and weak memory?

Effective Model Checking

Challenges:

- State Space Explosion

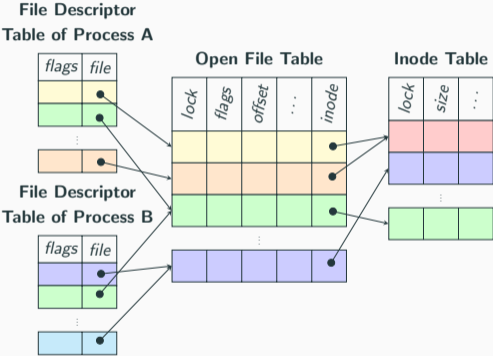
- **PerSeVerE**: framework for verification under ext4
 - formal model integrated with C/C++ consistency semantics
 - based on kernel's code and stress testing
 - effective model checking algorithm
- Report **bugs** in commonly used text editors like emacs, vim and nano

ext4 **Semantics**

Filesystem operations

Opening a file:

```
- df = open("foo.txt", O_APPEND)
```



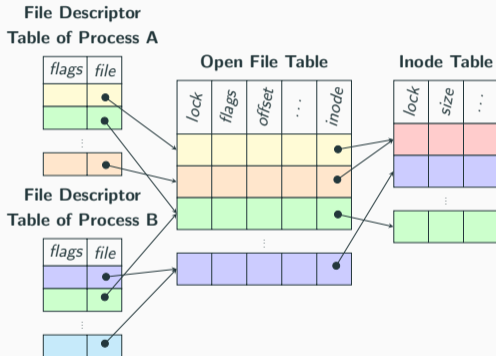
Filesystem operations

Opening a file:

- $d_f = \text{open}(\text{"foo.txt"}, \text{O_APPEND})$

Reading from a file:

- $r = \text{read}(d_f, 3), r = \text{pread}(d_f, 3, 42)$



Filesystem operations

Opening a file:

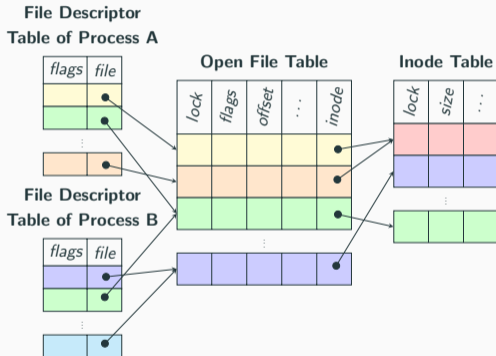
- $d_f = \text{open}(\text{"foo.txt"}, \text{O_APPEND})$

Reading from a file:

- $r = \text{read}(d_f, 3)$, $r = \text{pread}(d_f, 3, 42)$

Writing to a file:

- $\text{write}(d_f, \text{"foo"})$, $\text{pwrite}(d_f, \text{"bar"}, 0)$



Filesystem operations

Opening a file:

- $d_f = \text{open}(\text{"foo.txt"}, \text{O_APPEND})$

Reading from a file:

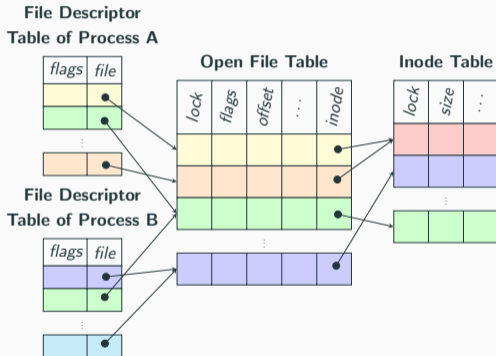
- $r = \text{read}(d_f, 3)$, $r = \text{pread}(d_f, 3, 42)$

Writing to a file:

- $\text{write}(d_f, \text{"foo"})$, $\text{pwrite}(d_f, \text{"bar"}, 0)$

Closing a file:

- $\text{close}(d_f)$



Filesystem operations

Opening a file:

- $d_f = \text{open}(\text{"foo.txt"}, \text{O_APPEND})$

Reading from a file:

- $r = \text{read}(d_f, 3), r = \text{pread}(d_f, 3, 42)$

Writing to a file:

- $\text{write}(d_f, \text{"foo"}), \text{pwrite}(d_f, \text{"bar"}, 0)$

Closing a file:

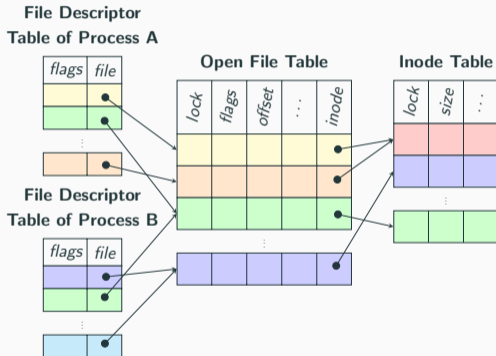
- $\text{close}(d_f)$

Directory operations:

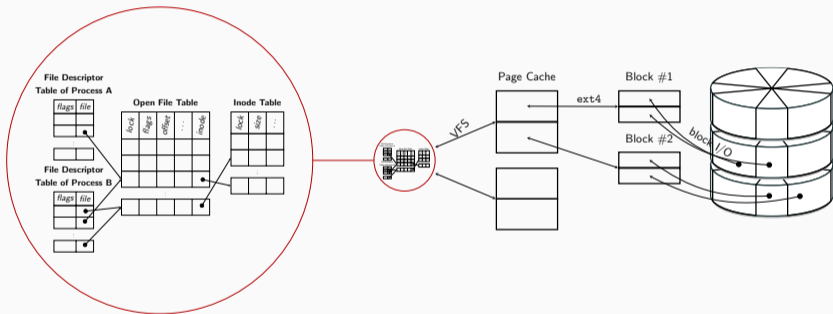
- $\text{creat}, \text{link}, \text{unlink}, \text{rename}$

Synchronization operations:

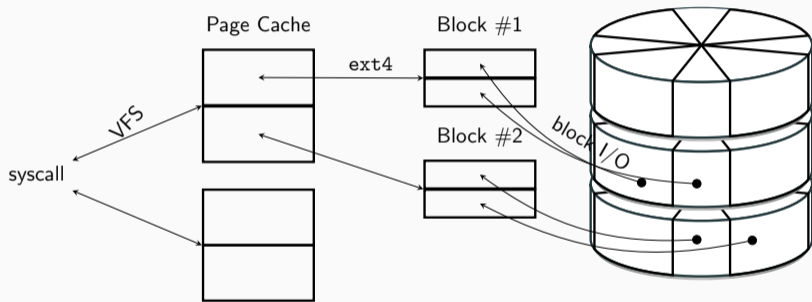
- $\text{fsync}, \text{sync}$



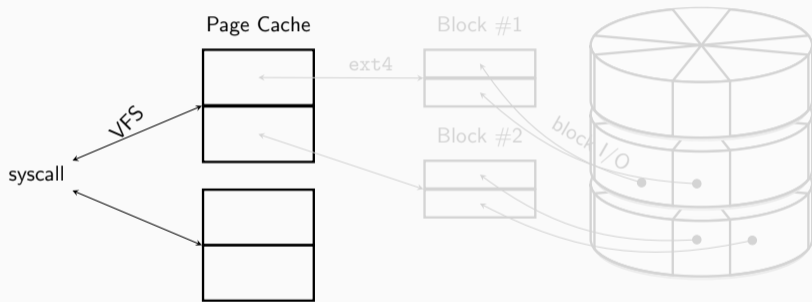
I/O Stack



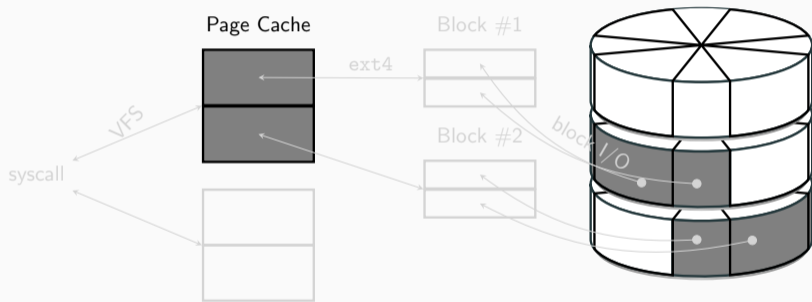
I/O Stack



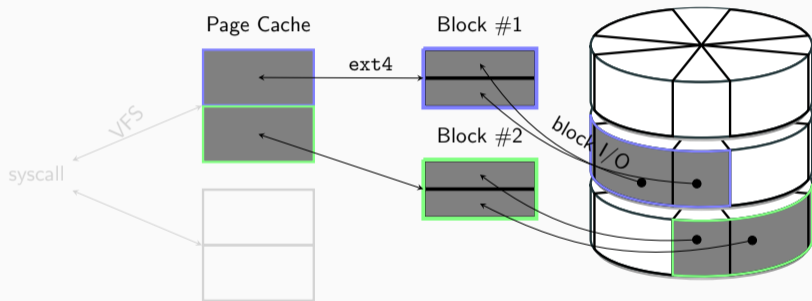
I/O Stack



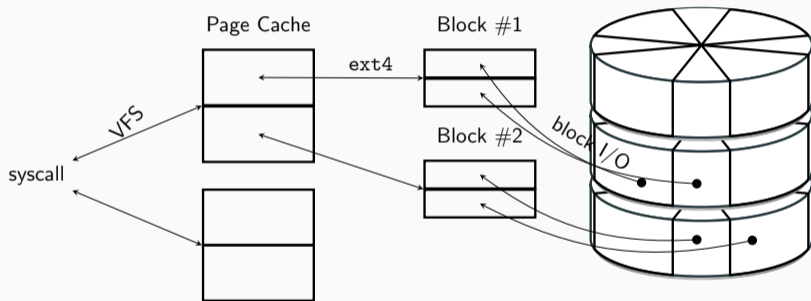
I/O Stack



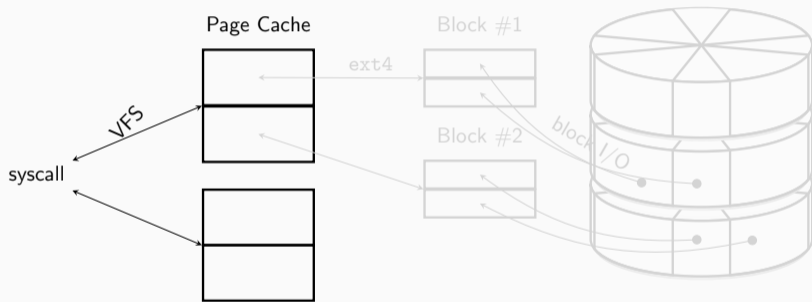
I/O Stack



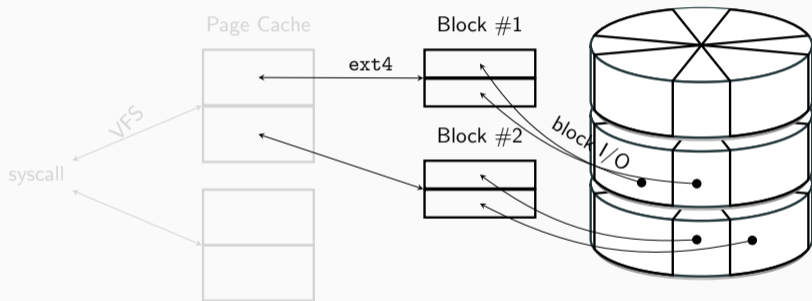
I/O Stack



I/O Stack



I/O Stack



Consistency Semantics

Writes vs Writes : c-atomic, c-ordered

$d_f \rightsquigarrow \text{"foo"}$

```
pwrite( $d_f$ , "bar", 0); || pwrite( $d_f$ , "qux", 0);  
r = pread( $d_f$ , 3, 0);
```

Possible outcomes: "bar", "qux"

Consistency Semantics

Writes vs Writes : c-atomic, c-ordered

Reads vs Overwrites :

```
 $d_f \rightsquigarrow \text{"foobar"}, \text{offset} = 0$ 
```

```
write( $d_f$ , "bar"); || r = read( $d_f$ , 3);
```

Possible outcomes: "foo", "bar"

Consistency Semantics

Writes vs Writes : c-atomic, c-ordered

Reads vs Overwrites :

$d_f \rightsquigarrow \text{"foo"}$

`write(df, "bar");` || `r = pread(df, 3, 0);`

Possible outcomes: "foo", "bar", "far", "fao", etc

Consistency Semantics

Writes vs Writes : c-atomic, c-ordered

Reads vs Overwrites : c-nonatomic, c-unordered

$d_f \rightsquigarrow$ "foo"

`write(d_f, "bar");` || `r = pread(d_f, 3, 0);`

Possible outcomes: "foo", "bar", "far", "fao", etc

Consistency Semantics

Writes vs Writes : c-atomic, c-ordered
Reads vs Overwrites : c-nonatomic, c-unordered
Reads vs Appends :

```
 $d_f \rightsquigarrow$  "foo", 0_APPEND, page_size = 3
```

```
write( $d_f$ , "barqux"); || r = pread( $d_f$ , 42, 0);
```

Possible outcomes: "foo", "foobar", "foobarqux"

Consistency Semantics

Writes vs Writes : c-atomic, c-ordered
Reads vs Overwrites : c-nonatomic, c-unordered
Reads vs Appends : page-c-atomic, page-c-ordered

```
 $d_f \rightsquigarrow \text{"foo"}, 0\_APPEND, page\_size = 3$ 
```

```
write( $d_f$ , "barqux"); || r = pread( $d_f$ , 42, 0);
```

Possible outcomes: "foo", "foobar", "foobarqux"

Consistency Semantics

Writes vs Writes : c-atomic, c-ordered
Reads vs Overwrites : c-nonatomic, c-unordered
Reads vs Appends : page-c-atomic, page-c-ordered
Directory vs All :

“foo.txt” \mapsto “foo”

```
db = creat (“foo.tmp”);  
write (db, “bar”); close (db);  
rename (“foo.tmp”, “foo.txt”);
```

```
df = open (“foo.txt”, O_RDONLY);  
r = read (df, 3);
```

Possible outcomes: “foo”, “bar”

Consistency Semantics

Writes vs Writes	:	c-atomic, c-ordered
Reads vs Overwrites	:	c-nonatomic, c-unordered
Reads vs Appends	:	page-c-atomic, page-c-ordered
Directory vs All	:	c-atomic, c-ordered

“foo.txt” \mapsto “foo”

```
db = creat (“foo.tmp”);  
write (db, “bar”); close (db);  
rename (“foo.tmp”, “foo.txt”);
```

```
df = open (“foo.txt”, O_RDONLY);  
r = read (df, 3);
```

Possible outcomes: “foo”, “bar”

Persistence Semantics

Overwrites :

```
 $d_f \rightsquigarrow$  "foo", sector_size = 1, block_size = 3
```

```
pwrite( $d_f$ , "bar", 0);
```

Possible outcomes: "foo", "boo", "bao", "bar" (not "fao", "far")

Persistence Semantics

Overwrites :

```
 $d_f \rightsquigarrow$  "foo", sector_size = 1, block_size = 1
```

```
pwrite( $d_f$ , "bar", 0);
```

Possible outcomes: "foo", "boo", "bao", "bar" (also "fao", "far")

Persistence Semantics

Overwrites : sector p-atomic, p-unordered

```
df  $\rightsquigarrow$  "foo", sector_size = 1, block_size = 1
```

```
pwrite(df, "bar", 0);
```

Possible outcomes: "foo", "boo", "bao", "bar" (also "fao", "far")

Persistence Semantics

Overwrites : sector p-atomic, p-unordered

Appends :

```
 $d_f \rightsquigarrow$  "foo", 0_APPEND, block_size = 3
```

```
pwrite( $d_f$ , "bar", 0);
```

Possible outcomes: "foo", "foobar"

Persistency Semantics

Overwrites : sector p-atomic, p-unordered

Appends : prefix-p-atomic *,

```
 $d_f \rightsquigarrow$  "foo", O_APPEND, block_size = 1
```

```
pwrite( $d_f$ , "bar", 0);
```

Possible outcomes: "foo", "foob", "fooba", "foobar"

Persistency Semantics

Overwrites : sector p-atomic, p-unordered

Appends : prefix-p-atomic *,

```
 $d_f \rightsquigarrow$  "foo", O_APPEND, block_size = 2
```

```
pwrite( $d_f$ , "bar", 0);
```

Possible outcomes: "foo", "foob", ~~"fooba"~~, "foobar" (also "foo0")

Persistence Semantics

Overwrites : sector p-atomic, p-unordered

Appends : prefix-p-atomic *,

```
 $d_f \rightsquigarrow$  "foo", O_APPEND, block_size = 2
```

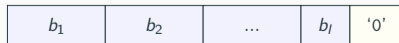
```
pwrite( $d_f$ , "bar", 0);
```

Possible outcomes: "foo", "foob", ~~"fooba"~~, "foobar" (also "foo0")

Before crash:



After crash:



Persistency Semantics

Overwrites : sector p-atomic, p-unordered

Appends : prefix-p-atomic *, same-file-p-ordered

```
 $d_a \rightsquigarrow \emptyset, d_b \rightsquigarrow \emptyset$ 
```

```
write( $d_a$ , "foo");
```

```
write( $d_b$ , "bar");
```

One possible outcome: $d_a \rightsquigarrow \emptyset, d_b \rightsquigarrow$ "bar"

Persistency Semantics

Overwrites : sector p-atomic, p-unordered
Appends : prefix-p-atomic *, same-file-p-ordered
{f}sync

$d_a \rightsquigarrow \emptyset, d_b \rightsquigarrow \emptyset$

`write(d_a , "foo");`

`fsync(d_a);`

`write(d_b , "bar");`

An impossible outcome: $d_a \rightsquigarrow \emptyset, d_b \rightsquigarrow$ "bar"

Persistency Semantics

Overwrites : sector p-atomic, p-unordered
Appends : prefix-p-atomic *, same-file-p-ordered
{f}sync : p-nonatomic, p-ordered

$d_a \rightsquigarrow \emptyset, d_b \rightsquigarrow \emptyset$

`write(d_a , "foo");`

`fsync(d_a);`

`write(d_b , "bar");`

An impossible outcome: $d_a \rightsquigarrow \emptyset, d_b \rightsquigarrow$ "bar"

Full model

- Axiomatic model in the style of RC11
- Assumes the consistency model contains an **hb** relation
- Includes a **pb** relation denoting the order in which disk accesses persist

Full model

- Axiomatic model in the style of RC11
- Assumes the consistency model contains an **hb** relation
- Includes a **pb** relation denoting the order in which disk accesses persist

$$\begin{aligned} & (I \cap D) \times (D \setminus I) \subseteq \mathbf{pb} && \text{(PB-INIT)} \\ & [DW]; (\mathbf{hb} \cap \mathbf{ssec}); [DW] \subseteq \mathbf{pb} && \text{(PB-SECTOR)} \\ & [DW]; (\mathbf{hb} \cap \mathbf{bseq}); [DW] \subseteq \mathbf{pb} && \text{(PB-BLOCK)} \\ & [DW_{\text{Floc}}]; (\mathbf{hb} \cap \mathbf{sf}); [DW_{\text{DsizeLoc}}] \subseteq \mathbf{pb} && \text{(PB-META)} \\ & [S \cup \text{FS}]; \mathbf{hb}; [D] \cup [D]; \mathbf{hb}; [S] \cup [DW]; (\mathbf{hb} \cap \mathbf{sf}); [\text{FS}] \subseteq \mathbf{pb} && \text{(PB-SYNC)} \\ & [DW_{\text{DnameLoc}} \cup DW^{\text{trunc}}]; \mathbf{hb}; [D \setminus DW_{\text{Floc}}] \subseteq \mathbf{pb} && \text{(PB-DIROPS)} \\ & (\mathbf{atom}; \mathbf{pb}) \cup (\mathbf{pb}; \mathbf{atom}) \subseteq \mathbf{pb} && \text{(PB-ATOM)} \end{aligned}$$

where

$$\mathbf{atom} \triangleq ([DW \setminus DW^{\text{zero}}]; (\mathbf{ssec} \cap \mathbf{sid}); [DW \setminus DW^{\text{zero}}]) \cup ([DW^{\text{rename}}]; \mathbf{sid}; [DW^{\text{rename}}]).$$

Full model

- Axiomatic model in the style of RC11
- Assumes the consistency model contains an **hb** relation
- Includes a **pb** relation denoting the order in which disk accesses persist

$$\begin{aligned} & (I \cap D) \times (D \setminus I) \subseteq \mathbf{pb} && \text{(PB-INIT)} \\ & [DW]; (\mathbf{hb} \cap \mathbf{ssec}); [DW] \subseteq \mathbf{pb} && \text{(PB-SECTOR)} \\ & [DW]; (\mathbf{hb} \cap \mathbf{bseq}); [DW] \subseteq \mathbf{pb} && \text{(PB-BLOCK)} \\ & [DW_{\text{Floc}}]; (\mathbf{hb} \cap \mathbf{sf}); [DW_{\text{DsizeLoc}}] \subseteq \mathbf{pb} && \text{(PB-META)} \\ & ([S \cup \text{FS}]; \mathbf{hb}; [D] \cup [D]; \mathbf{hb}; [S] \cup [DW]; (\mathbf{hb} \cap \mathbf{sf}); [\text{FS}] \subseteq \mathbf{pb} && \text{(PB-SYNC)} \\ & [DW_{\text{DnameLoc}} \cup DW^{\text{trunc}}]; \mathbf{hb}; [D \setminus DW_{\text{Floc}}] \subseteq \mathbf{pb} && \text{(PB-DIROPS)} \\ & (\mathbf{atom}; \mathbf{pb}) \cup (\mathbf{pb}; \mathbf{atom}) \subseteq \mathbf{pb} && \text{(PB-ATOM)} \end{aligned}$$

where

$$\mathbf{atom} \triangleq ([DW \setminus DW^{\text{zero}}]; (\mathbf{ssec} \cap \mathbf{sid}); [DW \setminus DW^{\text{zero}}]) \cup ([DW^{\text{rename}}]; \mathbf{sid}; [DW^{\text{rename}}]).$$

$$[S \cup \text{FS}]; \mathbf{hb}; [D]$$

Results

Save buffer as “foo.txt”

```
d_f = open (“foo.txt”, O_WRONLY|O_CREAT|O_TRUNC);  
write (d_f, buffer); fsync (d_f);  
close (d_f);  
printf (“File saved\n”);
```

The bug depends on the manifestation of a **race**

We **reported** the bug to the developers of nano

↪ proposed a fix and **verified** it with PerSeVerE

↪ our fixes were subsequently **merged**

We reproduced the same buggy pattern in emacs and vim

Summary

- **PerSeVerE**: framework for verification under ext4
 - formal model integrated with C/C++ consistency semantics
 - based on kernel's code and stress testing
 - effective model checking algorithm
- PerSeVerE is **available** at github.com/MPI-SWS/genmc

Future work

- Formalize other aspects of ext4
- Extend PerSeVerE for other filesystems