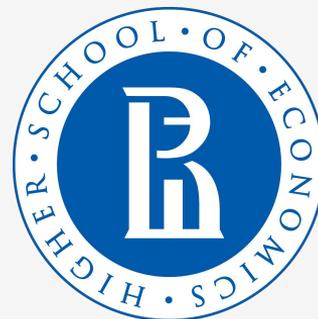# Persistency semantics of the ext4 filesystem

## Ilya Kaysin

**Research supervisor:**
**PhD Anton Podkopaev,**
**HSE, JetBrains Research**


**Reviewer:**
**PhD Christopher Pulte,**
**University of Cambridge**

NATIONAL RESEARCH
UNIVERSITY

**Memory**

```
// a = b = 0
a = 1
b = 1
// a = b = 1
```

## Memory

```
// a = b = 0
a = 1
b = 1
// a = b = 1
```

crash ⚡

```
// a = b = 0
```

## Memory

```
// a = b = 0
a = 1
b = 1
// a = b = 1
```

crash ⚡

```
// a = b = 0
```

## HDD

```
// a.txt = b.txt = "0"
a.txt <- "1"
b.txt <- "1"
// a.txt = b.txt = "1"
```

⚡

```
// a.txt = b.txt = "1"
```

**Memory**

```
// a = b = 0
a = 1
b = 1
// a = b = 1
```

crash

```
// a = b = 0
```

**HDD**

```
// a.txt = b.txt = "0"
a.txt <- "1"
b.txt <- "1"
// a.txt = b.txt = "1"
```

```
// a.txt = b.txt = "1"
//   "0"        "0"
//   "1"        "0"
```

**Memory**

```
// a = b = 0
a = 1
b = 1
// a = b = 1
```

crash⚡

```
// a = b = 0
```

**HDD**

```
// a.txt = b.txt = "0"
a.txt <- "1"
b.txt <- "1"
// a.txt = b.txt = "1"
```

⚡

```
// a.txt = b.txt = "1"
//    "0"        "0"
//    "1"        "0"
//    "0"        "1"
```

# Concurrent environment

```
// f.txt = 0

P
```

⚡

```
// f.txt = ???
```
⟵ - - - - - - - - - - -

**Persistency:** what is observable
after a system crash and restart

# Concurrent environment

```
// f.txt = 0

P1  ||  P2  || ... || Pn
```

// f.txt = **???**    ← ------------   **Persistency:** what is observable after a system crash and restart

# Concurrent environment

```
// f.txt = 0

P1  ||  P2  ||  ...  ||  Pn

// f.txt = ???  <------------------
```

**Consistency:** what is observable by concurrent threads

```
// f.txt = ???  <------------------
```

**Persistency:** what is observable after a system crash and restart

# Motivation

- Filesystems are complicated

- Ext4 – the most popular modern filesystem in Linux.

  Optimisations => obscure behavior => bugs

- For formal verification, **formal semantics** is required

# Motivation

- Filesystems are complicated

- Ext4 – the most popular modern filesystem in Linux.

  Optimisations => obscure behavior => bugs

- For formal verification, **formal semantics** is required

  mathematical model defining what can be observed (1) by other threads, (2) upon a crash

# Related work

**Weak persistency models**

- Persistency semantics of the Intel-x86 architecture [Raad et al. POPL 2019]

- Weak Persistency Semantics from the Ground Up [Raad et al. OOPSLA 2019]

**Filesystem models**

- Development of a Verified Flash File System [Schellhorn et al. ABZ 2014]

- Specifying and Checking File System Crash-Consistency Models [Bornholt et al. ASPLOS 2016]

# Related work

## Weak persistency models

- Persistency semantics of the Intel-x86 architecture [Raad et al. POPL 2019]

- Weak Persistency Semantics from the Ground Up [Raad et al. OOPSLA 2019]

## Filesystem models

- Development of a Verified Flash File System [Schellhorn et al. ABZ 2014]

- Specifying and Checking File System Crash-Consistency Models [Bornholt et al. ASPLOS 2016]

But no model encompasses **both:**

filesystem consistency and persistency

**Goal:** developing a formal model of ext4 compatible with formal verification

**Objectives:**

1. introduce a **general** filesystem **framework**

2. **specialize** the framework to build the **ext4 semantics**

3. adapt the ext4 semantics for **formal verification**

# General framework scheme:

**A program** → (1) → **Execution graphs** → (2) → **Consistent graphs** → (3) → **Persistent graphs**

# General framework scheme:

A program → (1) → **Execution** <u>**graphs**</u> → (2) → **Consistent graphs** → (3) → **Persistent graphs**
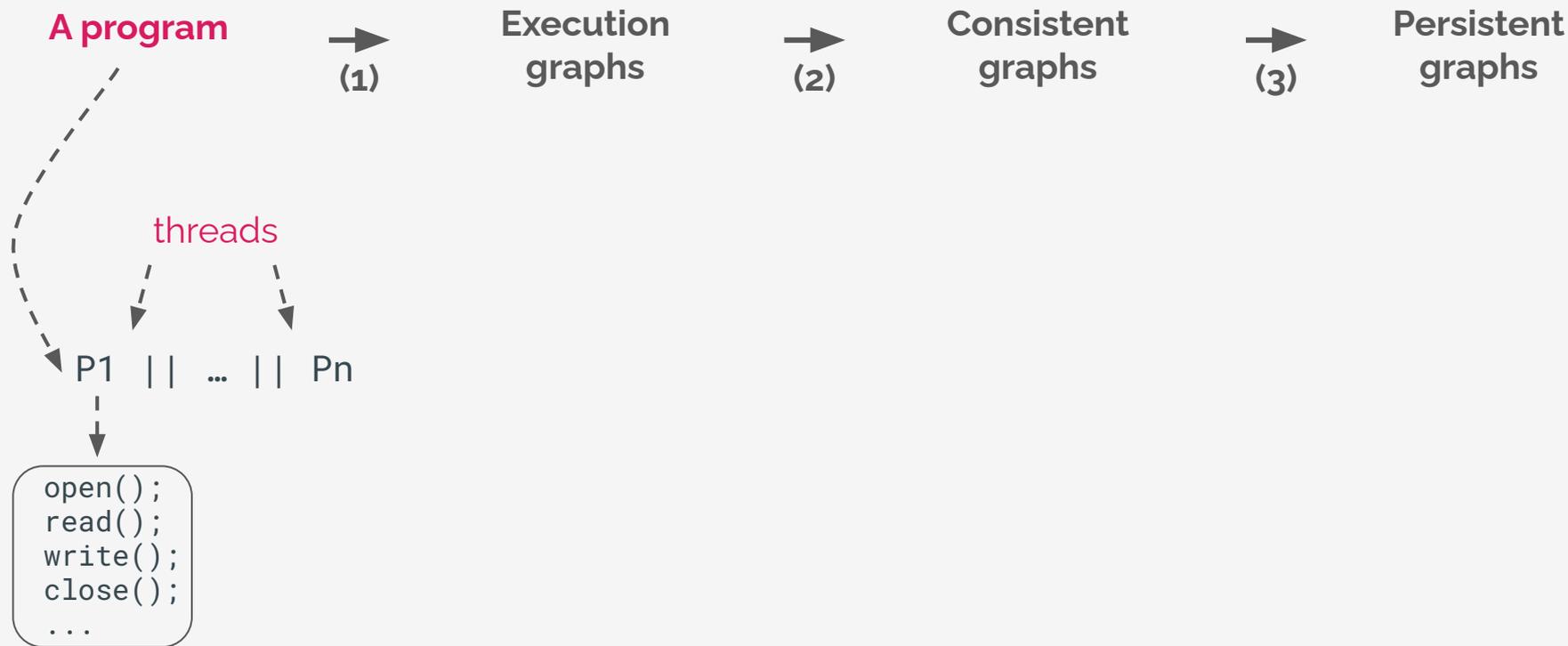
V = Events (Read, Write, ...)

E = Orders/Relations (**po**, **rf**, ...)

# General framework scheme:

**A program** → **(1)** → **Execution graphs** → **(2)** → **Consistent graphs** → **(3)** → **Persistent graphs**

threads

P1 || … || Pn

```
open();
read();
write();
close();
...
```

# General framework scheme:

**A program** → **(1)** **Execution graphs** → **(2)** **Consistent graphs** → **(3)** **Persistent graphs**

threads          recovery

P1 || ... || Pn  ⚡ REC

```
open();
read();
write();
close();
...
```

```
read();
read();
read();
read();
...
```

# General framework scheme:

**A program** → **(1)** → **Execution graphs** → **(2)** → **Consistent graphs** → **(3)** → **Persistent graphs**

threads     recovery

P1 || … || Pn ⚡ REC

```
open();
read();
write();
close();
...
```

```
read();
read();
read();
read();
...
```

**(2) Consistency**:

    Based on the C++ **memory** model

**(3) Persistency**:

- Defines **snapshots** - sets of writes surviving a crash

- Enforces the recovery to *read from a snapshot*

- Based on the *persists before* relation (**pb**) – the order in which writes become durable

## Objective 2

**Goal:** specializing the framework for ext4

1. Define **execution generation** from programs

2. Define **pb**

**Methodology:**

- consulting the **sources** (VFS + ext4)

- running **litmus tests**

## Objective 2

**Goal:** specializing the framework for ext4

1. Define **execution generation** from programs as a procedure generating *events* for each *system call*
2. Define **pb**

**Methodology:**

- consulting the **sources** (VFS + ext4)
- running **litmus tests**

## Objective 2

**Goal:** specializing the framework for ext4

1. Define **execution generation** from programs as a procedure generating *events* for each *system call*

2. Define **pb**

**Methodology:**

- consulting the **sources** (VFS + ext4)
- running **litmus tests**:
  - data race tests for consistency
  - write tests for persistency

||
V

Correctness

# Objective 2

**Goal:** specializing the framework for ext4

1. Define **execution generation** from programs as a procedure generating *events* for each *system call*

2. Define **pb** as the minimal relation satisfying the axioms:

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

**Methodology:**

- consulting the **sources** (VFS + ext4)

- running **litmus tests**:
  - data race tests for consistency
  - write tests for persistency

⇓

Correctness

# Objective 3

**Goal:** adapt ext4 semantics for effective formal verification

**(3) Persistency:**

**recovery reads from a snapshot**

Exponential

**Objective 3**

**Goal:** adapt ext4 semantics for effective formal verification

**(3) Persistency:**

**recovery reads from a snapshot**

Exponential

# Objective 3

**Goal:** adapt ext4 semantics for effective formal verification

**(3) Persistency:**

**recovery reads from a snapshot**

Exponential

**(3') Persistency:**

**executions do not contain certain paths**

$$[\mathrm{REC}]; \mathtt{rb}; \mathtt{atom}^?; \mathtt{pb}^?; \mathtt{rf}; [\mathrm{REC}] = \emptyset$$

# Objective 3

**Goal:** adapt ext4 semantics for effective formal verification

**(3) Persistency:**

**recovery reads from a snapshot**

Exponential

**(3') Persistency:**

**executions do not contain certain paths**

$$[\mathrm{REC}]; \mathtt{rb}; \mathtt{atom}^?; \mathtt{pb}^?; \mathtt{rf}; [\mathrm{REC}] = \emptyset$$

Equivalent to
(3) (proved)

Axiomatic

Can be used by
model checkers

# Results

- Built the general filesystem persistency+consistency framework

  (employs order **pb**; applicable to other filesystems and NVM)

- Constructed the ext4 semantics

  (by consulting the sources + running litmus tests)

- Adapted the semantics for formal verification

  (it helped to found anomalies in nano, vim, emacs)

- The work was published on POPL'21 (CORE A*)

M. Kokologiannakis, I. Kaysin, A. Raad, and V. Vafeiadis. **PerSeVerE: Persistency semantics for verification under ext4**. Proc. ACM Program. Lang., (POPL), 2021

# Backup slides

# Data loss anomalies in text editors

$save($ "f.txt" $) :$
    $d_f =$ open $($ "f.txt", O_FLAGS$)$;
    write $(d_f,$ BUF$)$;
    close $(d_f)$;

$backup($ "f.txt" $) :$
    $d_f =$ open $($ "f.txt", O_RDONLY$)$;
    $d_b =$ open $($ "f.txt$\sim$", O_FLAGS$)$;
    $b =$ read $(d_f)$;    write $(d_b, b)$;
    close $(d_b)$;    close $(d_f)$;

# **P**ersists **B**efore  (Ext4 specific)

$$(I \cap D) \times (D \setminus I) \subseteq pb \qquad \qquad (\textsc{pb-init})$$

$$[DW] ; (hb \cap ssec) ; [DW] \subseteq pb \qquad \qquad (\textsc{pb-sector})$$

$$[DW] ; (hb \cap bseq) ; [DW] \subseteq pb \qquad \qquad (\textsc{pb-block})$$

$$[DW_{Floc}] ; (hb \cap sf) ; [DW_{Dsizeloc}] \subseteq pb \qquad \qquad (\textsc{pb-meta})$$

$$[S \cup FS] ; hb ; [D] \ \cup \ [D] ; hb ; [S] \ \cup \ [DW] ; (hb \cap sf) ; [FS] \subseteq pb \qquad \qquad (\textsc{pb-sync})$$

$$[DW_{Dnameloc} \cup DW^{trunc}] ; hb ; [D \setminus DW_{Floc}] \subseteq pb \qquad \qquad (\textsc{pb-dirops})$$

$$(atom ; pb) \cup (pb ; atom) \subseteq pb \qquad \qquad (\textsc{pb-atom})$$

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

# System Call -> Events

write(a.txt, "hi") $\longrightarrow$

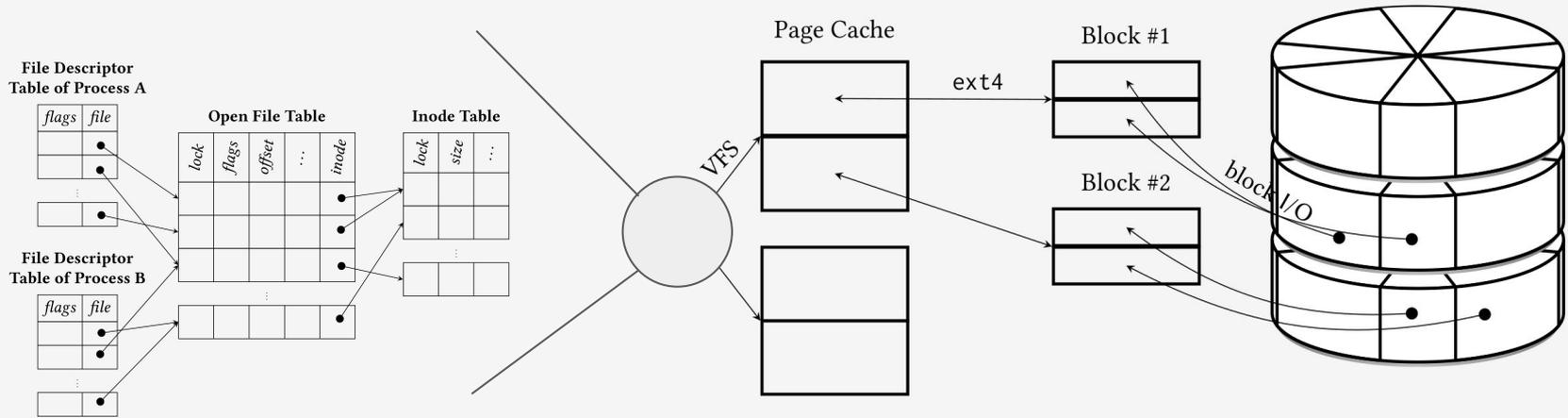$R_{size[a.txt]}$

$W_{a.txt[0]}$ "h"

$W_{a.txt[1]}$ "i"

$W_{size[a.txt]}$

# I/O Stack



File Descriptor
Table of Process A

| flags | file |
|-------|------|
| | ● |
| | ● |
| ⋮ | |
| | ● |

Open File Table

| lock | flags | offset | ⋯ | inode |
|------|-------|--------|---|-------|
| | | | | ● |
| | | | | ● |
| | | | | ● |
| ⋮ | | | | |
| | | | | ● |

Inode Table

| lock | size | ⋯ |
|------|------|---|
| | | |
| | | |
| ⋮ | | |
| | | |

File Descriptor
Table of Process B

| flags | file |
|-------|------|
| | ● |
| ⋮ | |
| | ● |

VFS

Page Cache

Block #1

ext4

Block #2

block I/O

# Consistent executions

1.  $\mathrm{cons}_{\mathcal{M}}(G)$

    $\mathrm{cons}_{\mathcal{M}}(.) \colon \mathrm{Exec} \to \{\mathrm{true}, \mathrm{false}\}$

    $\mathsf{hb} \triangleq (\mathsf{po} \cup \mathsf{sw})^{+}$ is irreflexive

    $(\mathsf{po} \cup \mathsf{rf})^{+}$ is irreflexive

2.  $G.\mathsf{pb}$ is irreflexive

**1: procedure** BufferRead($f$, *buf*, *count*, *o*)
2:  $\rightsquigarrow$ MR($ms_f$) **in** *size*
3:  $m \leftarrow \min(count, size - o)$
4:  **for** $i = 0$ **to** $m - 1$ **do**
5:   $\rightsquigarrow$ DR($(f, o+i)$) **in** *buf*[$i$]

**1: procedure** BufferWrite($f$, *buf*, *count*, *o*)
2:  $\rightsquigarrow$ L($f$)
3:  $\rightsquigarrow$ MR($ms_f$) **in** *size*
4:  **if** isPreallocBlock($o$, *count*, *size*) **then**
5:   $end \leftarrow \min(count + o, \texttt{getLastBlockEnd}(f))$
6:   **for** $i = size$ **to** $end - 1$ **do** $\rightsquigarrow$ DW$^{\texttt{zero}}$($(f, i), 0$)
7:   $\rightsquigarrow$ DW$^{\texttt{zero}}$($ds_f$, *end*)
8:   $size \leftarrow end$
9:  **if** $o > size$ **then**
10:   **for** $i = size$ **to** $o - 1$ **do** $\rightsquigarrow$ DW$^{\texttt{norm}}$($(f, i), 0$)
11:  **for** $i = 0$ **to** $count - 1$ **do**
12:   $\rightsquigarrow$ DW$^{\texttt{norm}}$($(f, o+i)$, *buf*[$i$])
13:   **if** (isFstB($f$, $o+i+1$) $\vee$ $i=count-1$)) $\wedge$ $o+i > size$ **then**
14:    $\rightsquigarrow$ DW$^{\texttt{norm}}$($ds_f$, $o+i$) ; MW($ms_f$, $o+i$)
15:  $\rightsquigarrow$ U($f$)

**1: procedure** pread($d_f$, *buf*, *count*, *o*)
2:  BufferRead($f$, *buf*, *count*, *o*)

**1: procedure** read($d_f$, *buf*, *count*)
2:  $\rightsquigarrow$ L($d_f$)
3:  $\rightsquigarrow$ MR($ol_{d_f}$) **in** $o$
4:  BufferRead($f$, *buf*, *count*, *o*)
5:  $\rightsquigarrow$ MW($ol_{d_f}$, $o + count$)
6:  $\rightsquigarrow$ U($d_f$)

**1: procedure** pwrite($d_f$, *buf*, *count*, *o*)
2:  BufferWrite($f$, *buf*, *count*, *o*)

**1: procedure** write($d_f$, *buf*, *count*)
2:  $\rightsquigarrow$ L($d_f$)
3:  $\rightsquigarrow$ MR($ol_{d_f}$) **in** $o$
4:  BufferWrite($f$, *buf*, *count*, *o*)
5:  $\rightsquigarrow$ MW($ol_{d_f}$, $o + count$)
6:  $\rightsquigarrow$ U($d_f$)

# System Calls

**Open:**
```
d = open("foo.txt", O_APPEND)
```

**Reading:**
```
r = read(d,3)
r = pread(d,3,42)
```

**Writing:**
```
write(d,"foo")
pwrite(d,"bar",0)
```

**Close:**
```
close(d)
```

**Synchronization:**
```
sync, fsync
```

**Directory operations:**
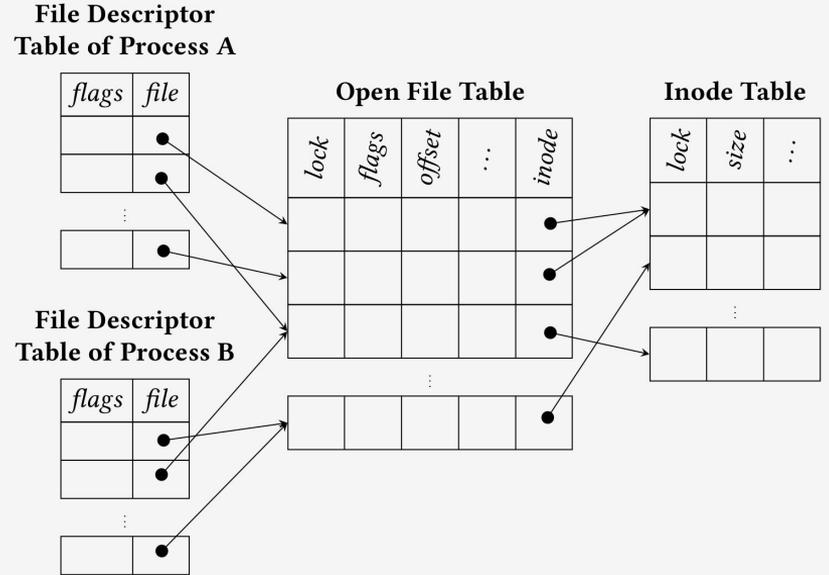```
create, link, unlink, rename
```

## Table 1. Sequential benchmarks used to evaluate PerSeVerE

| | P-Total | | P-Partial | | PerSeVerE | |
|---|---|---|---|---|---|---|
| | Execs | Time | Execs | Time | Execs | Time |
| pjnl-app-mfl | ⏱ | ⏱ | 15 | 57.39 | 9 | 0.02 |
| pjnl-crowr-ord | 1323 | 7.84 | 18 | 0.04 | 6 | 0.02 |
| pjnl-owr-at | 36 | 0.03 | 5 | 0.02 | 5 | 0.02 |
| pord-owr-N | ⏱ | ⏱ | 243 | 308.26 | 2 | 0.02 |
| pord-app-mbl | 12 | 0.03 | 5 | 0.03 | 3 | 0.02 |
| pord-app-sfl | 900 | 1.90 | 13 | 0.03 | 3 | 0.02 |
| pord-MP+fsync | 8724 | 496.63 | 3 | 0.03 | 3 | 0.02 |
| pord-trapp-dfl-ord | 4044 | 47.22 | 5 | 0.03 | 3 | 0.02 |
| pord-crapp-ooo | 810 | 1.38 | 21 | 0.03 | 3 | 0.02 |
| pord-MP+osync | 8724 | 596.20 | 3 | 0.03 | 3 | 0.02 |
| pord-owrapp-ord | 450 | 0.39 | 10 | 0.03 | 4 | 0.02 |
| pord-owrapp-ord2 | 3780 | 26.69 | 21 | 0.03 | 5 | 0.02 |
| pord-owr-sbl | 592 | 0.26 | 25 | 0.02 | 6 | 0.02 |
| pord-rnm-at2 | 696 | 0.14 | 9 | 0.06 | 5 | 0.02 |
| pord-rnmtrapp-ord | 3240 | 5.87 | 18 | 0.05 | 6 | 0.02 |
| pord-app-N | ⏱ | ⏱ | ⏱ | ⏱ | 5 | 0.03 |
| nano-backup-old | ⏱ | ⏱ | ⏱ | ⏱ | 46 | 0.08 |
| nano-backup-fix | ⏱ | ⏱ | ⏱ | ⏱ | 10 | 0.04 |

## Table 2. Concurrent benchmarks used to evaluate PerSeVerE

| | P-Total | | P-Partial | | PerSeVerE | |
|---|---|---|---|---|---|---|
| | Execs | Time | Execs | Time | Execs | Time |
| pord-wr+rdwr+fsync | ⏱ | ⏱ | 206 | 0.30 | 6 | 0.03 |
| pord-wr+wr-N | ⏱ | ⏱ | ⏱ | ⏱ | 6240 | 2.13 |
| pord-wr+wr-N-RR | ⏱ | ⏱ | ⏱ | ⏱ | 17 | 0.72 |
| pord-wr+wr-N-unord | ⏱ | ⏱ | ⏱ | ⏱ | 22 680 | 27.74 |
| pord-wr+wr-N-join-main | ⏱ | ⏱ | ⏱ | ⏱ | 216 | 1.19 |
| pord-wr+wr-N-join-thr | ⏱ | ⏱ | ⏱ | ⏱ | 2052 | 37.00 |
| pord-rd-wr+wr-N-cont | ⏱ | ⏱ | ⏱ | ⏱ | 12 888 | 61.99 |
| pord-rd-wr+wr-N-join | ⏱ | ⏱ | ⏱ | ⏱ | 216 | 3.76 |

# Пример: перезапись

```
// a.txt = "000...0"
a.txt <- "111...1"
```

⚡

```
          a.txt
"00000000000"
"110011110000"
"111100001100"
"1111111111111"
          ...
```

- **Сектора записываются атомарно**
- **Сектора внутри одного блока упорядочены**