technische universität
dortmund

# *Beastie In For Checkup: Analyzing FreeBSD with LockDoc*

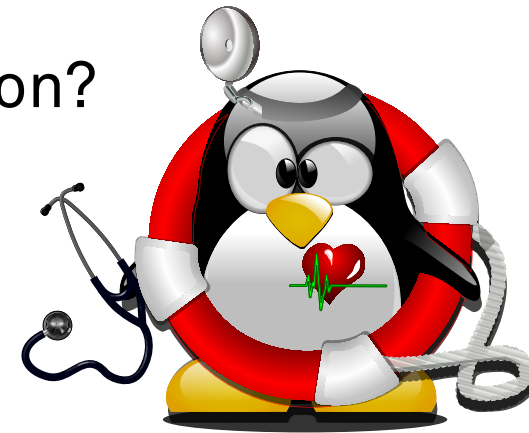**Alexander Lochmann,** Horst Schirmeier

alexander.lochmann@tu-dortmund.de
https://ess.cs.tu-dortmund.de/~al

Databases and Information Systems Group
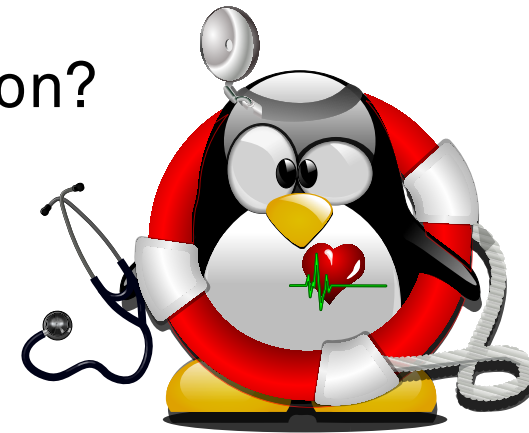Computer Science 6, TU Dortmund

# What is LockDoc?

- Tracks **locking pattern** and **data-structure** accesses

- **Recording** performed under a **load**

- Generates documentation, and locates locking bugs

- **Validate** existing locking **documentation**

  - Does the code adhere to the documentation?

- LockDoc study on Linux [3]

  - Validate documentation of 5 data type

  - 53 % of all observed fields accessed consistently with their doc.

Alexander Lochmann, Horst Schirmeier, Hendrik Borghorst, and Olaf Spinczyk. 2019. *LockDoc: Trace-Based Analysis of Locking in the Linux Kernel*. EuroSys'19.
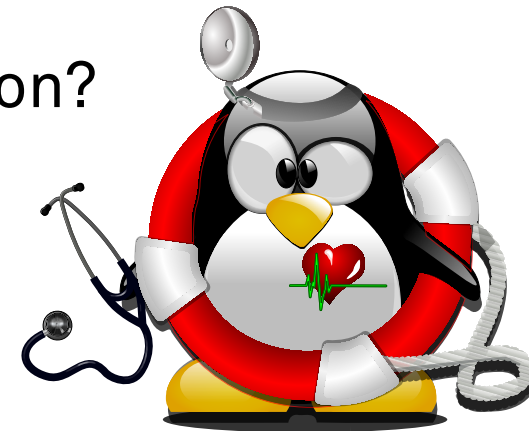
# What is LockDoc?

- Tracks **locking pattern** and **data-structure** accesses

- **Recording** performed under a **load**

- Generates documentation, and locates locking bugs

- **Validate** existing locking **documentation**

  - Does the code adhere to the documentation?

- LockDoc study on Linux [3]

  - Validate documentation of 5 data type

  - 53 % of all observed fields accessed consistently with their doc.

- Word-size variables can be accessed without locks
- If no concurrency takes place, no locks needed
- No locks if consistency does not matter

# What is LockDoc?

- Tracks **locking pattern** and **data-structure** accesses

- **Recording** performed under a **load**

- Generates documentation, and locates locking bugs

- **Validate** existing locking **documentation**

  - Does the code adhere to the documentation?

- LockDoc study on Linux [3]

  - Validate documentation of 5 data type

  - 53 % of all observed fields accessed consistently with their doc.

→ **Real bugs?**
→ **Issues with LockDoc?**

Alexander Lochmann, Horst Schirmeier, Hendrik Borghorst, and Olaf Spinczyk. 2019. *LockDoc: Trace-Based Analysis of Locking in the Linux Kernel*. EuroSys'19.

# Locking Documentation in FreeBSD

```
/*
 * Reading or writing any of these items requires
   holding the appropriate lock.
 *
 * Lock reference:
 *   c - namecache mutex
 *   i - interlock
 *   l - mp mnt_listmtx or freelist mutex
 *   I - updated with atomics, 0->1 and 1→0
         transitions with interlock held
 *   m - mount point interlock
 *   p - pollinfo lock
 *   u - Only a reference to the vnode is needed to
         read.
 *   v - vnode lock
 *
 * Vnodes may be found on many lists.  The general way
   to deal with operating
 * on a vnode that is on a list is:
 *   1) Lock the list and find the vnode.
 *   2) Lock interlock so that the vnode does not go
         away.
 *   3) Unlock the list to avoid lock order reversals.
 *   4) vget with LK_INTERLOCK and check for ENOENT, or
 *   5) Check for DOOMED if the vnode lock is not
         required.
 *   6) Perform your operation, then vput().
 */
```

sys/sys/vnode.h

# Locking Documentation in FreeBSD

```
/*
 * Reading or writing any of these items requires
   holding the appropriate lock.
 *
 * Lock reference:
 *   c - namecache mutex
 *   i - interlock
 *   l - mp mnt_listmtx or freelist mutex
 *   I - updated with atomics, 0->1 and 1→0
        transitions with interlock held
 *   m - mount point interlock
 *   p - pollinfo lock
 *   u - Only a reference to the vnode is needed to
        read.
 *   v - vnode lock
 *
 * Vnodes may be found on many lists.  The general way
   to deal with operating
 * on a vnode that is on a list is:
 *   1) Lock the list and find the vnode.
 *   2) Lock interlock so that the vnode does not go
        away.
 *   3) Unlock the list to avoid lock order reversals.
 *   4) vget with LK_INTERLOCK and check for ENOENT, or
 *   5) Check for DOOMED if the vnode lock is not
        required.
 *   6) Perform your operation, then vput().
 */
```

sys/sys/vnode.h

# Locking Documentation in FreeBSD

```
/*
 * Reading or writing any of these items requires
   holding the appropriate lock.
 *
 * Lock reference:
 *   c - namecache mutex
 *   i - interlock
 *   l - mp mnt_listmtx or freelist mutex
 *   I - updated with atomics, 0->1 and 1→0
         transitions with interlock held
 *   m - mount point interlock
 *   p - pollinfo lock
 *   u - Only a reference to the vnode is needed to
         read.
 *   v - vnode lock
 *
 * Vnodes may be found on many lists.  The general way
   to deal with operating
   on a vnode that is on a list is:
 *   1) Lock the list and find the vnode.
 *   2) Lock interlock so that the vnode does not go
      away.
 *   3) Unlock the list to avoid lock order reversals.
 *   4) vget with LK_INTERLOCK and check for ENOENT, or
 *   5) Check for DOOMED if the vnode lock is not
      required.
 *   6) Perform your operation, then vput().
 */
```

sys/sys/vnode.h

# Locking Documentation in FreeBSD

```
/*
 * Reading or writing any of these items requires
   holding the appropriate lock.
 *
 * Lock reference:
 *   c - namecache mutex
 *   i - interlock
 *   l - mp mnt_listmtx or freelist mutex
 *   I - updated with atomics, 0->1 and 1→0
         transitions with interlock held
 *   m - mount point interlock
 *   p - pollinfo lock
 *   u - Only a reference to the vnode is needed to
         read.
 *   v - vnode lock
 *
 * Vnodes may be found on many lists.  The general way
   to deal with operating
 * on a vnode that is on a list is:
 *   1) Lock the list and find the vnode.
 *   2) Lock interlock so that the vnode does not go
         away.
 *   3) Unlock the list to avoid lock order reversals.
 *   4) vget with LK_INTERLOCK and check for ENOENT, or
 *   5) Check for DOOMED if the vnode lock is not
         required.
 *   6) Perform your operation, then vput().
 */
```
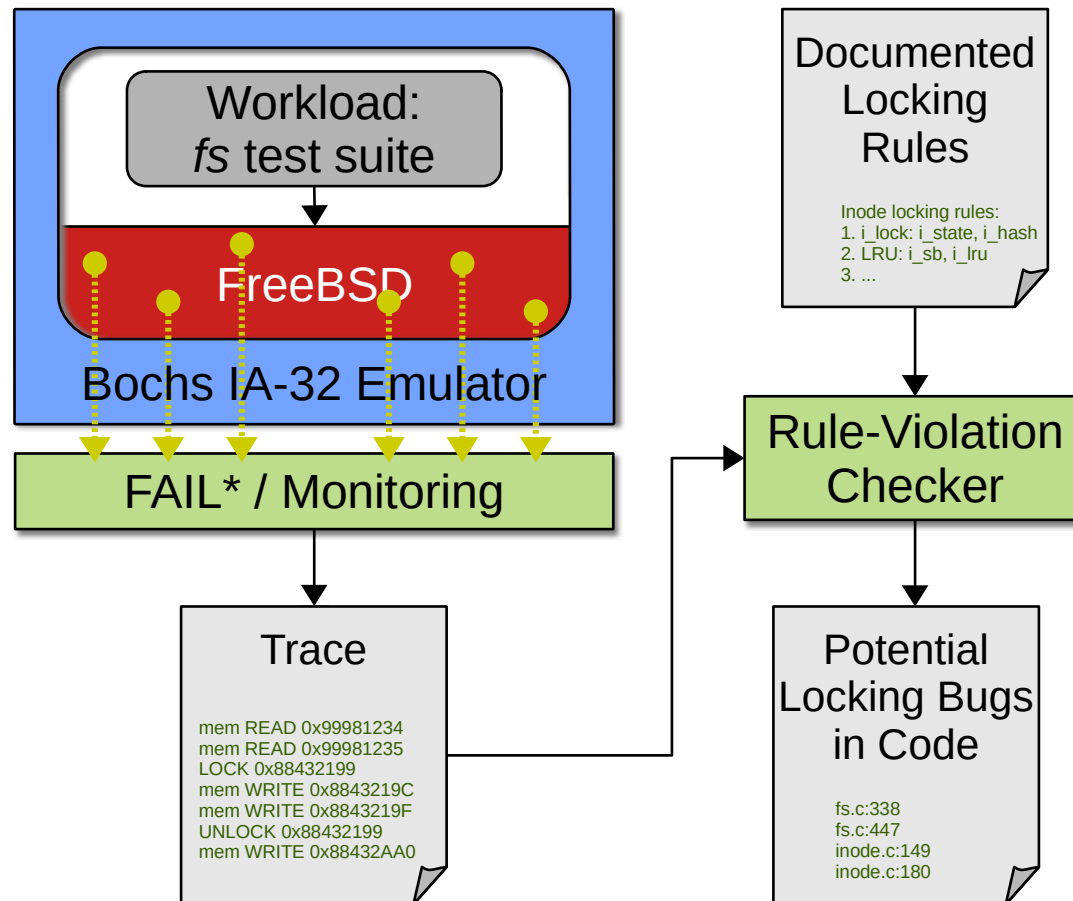
sys/sys/vnode.h

# Locking Documentation in FreeBSD

```
/*
 * Reading or writing any of these items requires
   holding the appropriate lock.
 *
 * Lock reference:
 *   c - namecache mutex
 *   i - interlock
 *   l - mp mnt_listmtx or freelist mutex
 *   I - updated with atomics, 0->1 and 1→0
          transitions with interlock held
 *   m - mount point interlock
 *   p - pollinfo lock
 *   u - Only a reference to the vnode is needed to
          read.
 *   v - vnode lock
 *
 * Vnodes may be found on many lists.  The general way
   to deal with operating
 * on a vnode that is on a list is:
 *   1) Lock the list and find the vnode.
 *   2) Lock interlock so that the vnode does not go
       away.
 *   3) Unlock the list to avoid lock order reversals.
 *   4) vget with LK_INTERLOCK and check for ENOENT, or
 *   5) Check for DOOMED if the vnode lock is not
       required.
 *   6) Perform your operation, then vput().
 */
```

sys/sys/vnode.h

```c
struct vnode {
    // […]
    short   v_irflag;/* i frequently read flags */
    seqc_t  v_seqc;  /* i modification count */
    // […]
    /*
     * Filesystem instance stuff
     */
    struct mount *v_mount; /* u [...] */
    TAILQ_ENTRY(vnode) v_nmntvnodes; /* m [...] */
    // [...]
};
```

# Locking Documentation in FreeBSD

```
/*
 * Reading or writing any of these items requires
   holding the appropriate lock.
 *
 * Lock reference:
 *   c - namecache mutex
 *   i - interlock
 *   l - mp mnt_listmtx or freelist mutex
 *   I - updated with atomics, 0->1 and 1→0
        transitions with interlock held
 *   m - mount point interlock
 *   p - pollinfo lock
 *   u - Only a reference to the vnode is needed to
        read.
 *   v - vnode lock
 *
 * Vnodes may be found on many lists.  The general way
   to deal with operating
 * on a vnode that is on a list is:
 *   1) Lock the list and find the vnode.
 *   2) Lock interlock so that the vnode does not go
        away.
 *   3) Unlock the list to avoid lock order reversals.
 *   4) vget with LK_INTERLOCK and check for ENOENT, or
 *   5) Check for DOOMED if the vnode lock is not
        required.
 *   6) Perform your operation, then vput().
 */
```

sys/sys/vnode.h

```
struct vnode {
    // [...]
    short    v_irflag;/* i frequently read flags */
    seqc_t   v_seqc;  /* i modification count */
    // [...]
    /*
     * Filesystem instance stuff
     */
    struct mount *v_mount; /* u [...] */
    TAILQ_ENTRY(vnode) v_nmntvnodes; /* m [...] */
    // [...]
};
```

# LockDoc – A Different Approach

# Instrumentation / Experiment Setup

- i386 FreeBSD 13.0 (Git commit *2134e85bc*)

- Instrument FreeBSD's *Witness* system [2]

  – Uses same lock model as LockDoc: read lock, write lock, rw lock

  – Automatically instrument all lock operations

  – 8 different types recorded: *hardirq*, *lockmgr*, *rm*, *rw*, *sleepable rm*, *sleep mutex*, *spin mutex*, and *sx*

- Using *fs* test suite from *Linux Test Project* as workload

- 26.43 hours runtime (20.22 minutes in a real vm)

# FreeBSD Results (1)

- 4 data types: *vnode*, *mount*, *buf*, and *bufobj*

- *bufobj* is embedded in *vnode*

| Data Type | #R | #No | #Ob | ✓ (%) | ~ (%) | ✗ (%) |
|---|---|---|---|---|---|---|
| vnode | 82 | 9 | 73 | 72.60 | 27.40 | 0.00 |
| mount | 38 | 7 | 31 | 74.19 | 25.81 | 0.00 |
| buf | 80 | 10 | 70 | 71.43 | 27.14 | 1.43 |

# FreeBSD Results (1)

- 4 data types: *vnode*, *mount*, *buf*, and *bufobj*
- *bufobj* is embedded in *vnode*

| Data Type | #R | #No | #Ob | ✓ (%) | ~ (%) | ✗ (%) |
|---|---|---|---|---|---|---|
| vnode | 82 | 9 | 73 | 72.60 | 27.40 | 0.00 |
| mount | 38 | 7 | 31 | 74.19 | 25.81 | 0.00 |
| buf | 80 | 10 | 70 | 71.43 | 27.14 | 1.43 |

Read of *buf.b_error*

# FreeBSD Results (1)

- 4 data types: *vnode*, *mount*, *buf*, and *bufobj*
- *bufobj* is embedded in *vnode*

| Data Type | #R | #No | #Ob | ✓ (%) | ~ (%) | ✗ (%) |
|-----------|----|-----|-----|-------|-------|-------|
| vnode | 82 | 9 | 73 | 72.60 | 27.40 | 0.00 |
| mount | 38 | 7 | 31 | 74.19 | 25.81 | 0.00 |
| buf | 80 | 10 | 70 | 71.43 | 27.14 | 1.43 |

72.41 %

Beastie In For Checkup: Analyzing FreeBSD with LockDoc

# FreeBSD Results (1)

- 4 data types: *vnode*, *mount*, *buf*, and *bufobj*
- *bufobj* is embedded in *vnode*

| Data Type | #R | #No | #Ob | ✓ (%) | ~ (%) | ✗ (%) |
|-----------|-----|-----|-----|--------|--------|--------|
| vnode | 82 | 9 | 73 | 72.60 | 27.40 | 0.00 |
| mo | | | | | | .00 |
| buf | | | | | | .43 |

→ **What about the remaining 27.59 %?**

# FreeBSD Results (2)

$$0{,}9 \leq s_r < 1$$

[1]https://lists.freebsd.org/archives/freebsd-fs/2021-August/000371.html
[2]https://github.com/freebsd/freebsd-src/blob/main/sys/ufs/ufs/inode.h#L75

# FreeBSD Results (2)

- Inspecting tuples with relative support $0,9 \leq s_r < 1$

    - 11 tuples found

    - 9 false positives

        - No locks needed due to domain-specific knowledge[1]

        - Unguarded NULL-pointer checks

        - Locking pattern not covered by LockDoc[2]:

            a Acquire *vnode lock* exclusively
            b Use *vnode lock* in shared mode **+** *interlock*

[1]https://lists.freebsd.org/archives/freebsd-fs/2021-August/000371.html
[2]https://github.com/freebsd/freebsd-src/blob/main/sys/ufs/ufs/inode.h#L75

# FreeBSD Results (2)

- Inspecting tuples with relative support $0,9 \leq s_r < 1$

  - 11 tuples found

  - 9 false positives

    - No locks needed due to domain-specific knowledge[1]

    - Unguarded NULL-pointer checks

    - Locking pattern not covered by LockDoc[2]:
      - a  Acquire *vnode lock* exclusively
      - b  Use *vnode lock* in shared mode **+** *interlock*

  - 2  real bugs

    - Relative support of 97.3 % and 96.2 %

    - Unguarded write to *buf.b_vflags* and read of *buf.b_blkno*

[1]https://lists.freebsd.org/archives/freebsd-fs/2021-August/000371.html
[2]https://github.com/freebsd/freebsd-src/blob/main/sys/ufs/ufs/inode.h#L75

# Locking Bug 1

- Unguarded write to *buf.b_vflags*



https://github.com/freebsd/freebsd-src/commit/e3d675958539eee899d42438f5b46a26f3c64902

# Locking Bug 2

- Unguarded read of *buf.b_blkno*

# Locking Bug 2

- Unguarded read of *buf*.*b_blkno*

# FreeBSD Results (2)

- Inspecting tuples with relative support $0,9 \leq s_r < 1$

  - 11 tuples found

  - 9 false positives

    - No locks needed due to domain-specific knowledge

    - Unguarded NULL-pointer checks

    - Pattern not covered by LockDoc:

      a  Acquire *vnode lock* exclusively
      b  Use *vnode lock* in shared mode **+** *interlock*

  - 2  real bugs

    - Relative support of 97.3 % and 96.2 %

    - Unguarded write to *buf.b_vflags* and read of *buf.b_blkno*

- *Taking samples from  tuples with relative support* $s_r < 0.9$

  - *buf.b_qindex and buf.b_subqueue are "Protected by the buf queue lock"[1]*

[1] cf. sys/sys/buf.h, line 96

# FreeBSD Results (3)

- *buf.b_qindex* and *buf.b_subqueue are*

  *"Protected by the buf queue lock"[1]*

The "buf queue lock":
```
struct bufqueue {
    struct mtx_padalign   bq_lock;
    // [...]
};
```

Multiple buf queues exist:
```
struct bufdomain {
    struct bufqueue       bd_subq[MAXCPU + 1]; /* [...] */
    struct bufqueue       bd_dirtyq;
    struct bufqueue      *bd_cleanq;
    struct mtx_padalign bd_run_lock;
    // [...]
};
```

- Multiple locks exist:
  - *bq_subq.bq_lock*
  - *bq_dirtyq.bq_lock*
- Accesses are split across them

[1] cf. sys/sys/buf.h, line 96

# Summary

- 72.41 % of all observed fields adhere to locking documentation

- Using sound locking documentation to search for bugs

    – Found limitations of LockDoc

    – Found **2 locking bugs**

- Outlook

    – Integrate lock classes, e.g., *bq_subq.bq_lock ↔ bq_dirtyq.bq_lock*

    – Further investigate rules with rel. support < 90 %

# References

[1] Robert Love. 2010. *Linux Kernel Development (3rd ed.)*.

[2] Marshall Kirk McKusick, George V. Neville-Neil, and Robert N. M. Watson. 2014. *The Design and Implementation of the FreeBSD Operating System*.

[3] Alexander Lochmann, Horst Schirmeier, Hendrik Borghorst, and Olaf Spinczyk. 2019. *LockDoc: Trace-Based Analysis of Locking in the Linux Kernel*. EuroSys'19.

[4] https://github.com/linux-test-project/ltp