

Nested Locks in the Lock Implementation: The Real-Time Read-Write Semaphores on Linux

Daniel B. de Oliveira^{1,2,3}, Daniel Casini³, Rômulo S. de Oliveira², Tommaso Cucinotta³,
Alessandro Biondi³, and Giorgio Buttazzo³

¹RHEL Platform/Real-time Team, Red Hat, Inc., Pisa, Italy.

²Department of Systems Automation, UFSC, Florianópolis, Brazil.

³RETIS Lab, Scuola Superiore Sant'Anna, Pisa, Italy.

Email: bristol@redhat.com, romulo.deoliveira@ufsc.br,

{daniel.casini,tommaso.cucinotta, alessandro.biondi, giorgio.buttazzo}@santannapisa.it

I. INTRODUCTION

Linux is a general purpose operating system (GPOS) that has gained many real-time (RT) features over the last decade. For instance, nowadays Linux has a *fully preemptive* mode and a deadline-oriented scheduler [1]. Although some of these features are part of the official Linux kernel, many of them are still part of an external patch set, the PREEMPT-RT [2]. The PREEMPT-RT changes the locking methods of Linux to prevent unbounded priority inversion. This is achieved by using the Priority Inheritance Protocol [3] on in-kernel mutexes, which bounds the activation delay in high priority tasks. Indeed, the *latency* is the main evaluation metric for the PREEMPT-RT Linux: for example, the Red Hat Enterprise Linux for Real-time [4] (based on PREEMPT-RT) shows a maximum latency of 150 μ s on certified hardware. However, due to Linux's GPOS nature, RT Linux developers are challenged to provide the predictability required for an RTOS, while not causing regressions on the general purpose benchmarks. As a consequence, the implementation of some well known algorithms, like read/write semaphores, has been done using approaches that were not well explored in academic papers.

II. READ-WRITE SEMAPHORES ON LINUX

On Linux, the read-write semaphores provide concurrent readers and exclusive writers for a given critical section. For example, since the memory mapping information of a process is read very often but rarely changes during its execution, it is protected by a read-write semaphore.

The API of the read-write semaphores is composed by four main functions. Readers call `DOWN_READ()` before entering in the read-side, calling `UP_READ()` when leaving the read-side of the critical section. Writers should call `DOWN_WRITE()` before entering in the write-side of the critical section, calling `UP_WRITE()` when leaving. These functions take only one argument, which is a pointer to a structure `rw_semaphore`. The `rw_semaphore` structure is presented in Figure 1¹.

The `readers` variable is an atomic type that counts how many concurrent readers are inside the critical section. This variable is also used to store `READER BIAS` and `WRITER BIAS` flags, which are used to define if there are either readers or a writer in the critical section. Whenever a task should block in the semaphore, it will do by blocking in the real-time mutex `rt_mutex` of the semaphore. The `rt_mutex` is defined as shown in Figure 2¹:

```
1 struct rw_semaphore {
2   atomic_t          readers;
3   struct rt_mutex   rtmutex;
4};
```

Fig. 1: Read-write Semaphore structure

```
1 struct rt_mutex {
2   raw_spinlock_t   wait_lock;
3   struct rb_root_cached waiters;
4   struct task_struct *owner;
5   int               save_state;
6};
```

Fig. 2: Real-time Mutex structure

In order to protect the fields of the `rt_mutex` struct from concurrent accesses, the spin lock `wait_lock` is used whenever the internal fields of the mutex are modified. The `wait_lock` of the real-time mutex is also used to avoid two writers setting the `WRITE/READ BIAS` concurrently in the `rw_semaphore` structure.

The pseudo-code of each operation is presented in Figure 3 and 4, respectively.

III. OPEN PROBLEMS

Considering our example, when `DOWN_WRITE()` is called, the task that is trying to acquire the read/write semaphore for writing has to lock two nested resources, a regular *mutex* (acquired at line 8, Figure 4) and a *spin lock* (acquired at line 14, Figure 4), thus creating a **heterogeneous nested lock** (e.g., a suspension-based lock with a nested spin-based lock or vice-versa). This case study, taken from the Linux kernel, highlights two open issues. The first one concerns the need for implementing

¹Debug fields removed from structure's definition.

```

1: function UP_READ(rw_sem) /* using atomic operations */
2:   if --rw_sem->readers == 0 then
3:     if a writer is holding the rw_sem->rtmutex then
4:       wake-up the writer
5:     end if
6:   end if
7: end function
8:
9: function DOWN_READ(rw_sem)
10:  if ++rw_sem->readers > 1 then /* using atomic operations */
11:    return /* enter the critical section */
12:  else
13:    rw_sem->readers--
14:  end if
15:  take rw_sem->rtmutex.wait_lock /* might block busy (spinlock) */
16:  if WRITER BIAS is not set then
17:    rw_sem->readers++
18:    release rw_sem->rtmutex.wait_lock
19:    return /* enter in the critical section */
20:  end if
21:  release rw_sem->rtmutex.wait_lock
22:  take rw_sem->rt_mutex /* might block suspended (rt mutex) */
23:  rw_sem->readers++
24:  release the rw_sem->rt_mutex
25:  return /* enter in the critical section */
26: end function

```

Fig. 3: Read-side operations

```

1: function UP_WRITE(rw_sem)
2:   clear WRITER BIAS
3:   set READER BIAS
4:   release sem->rtmutex
5: end function
6:
7: function DOWN_WRITE(rw_sem)
8:   take rw_sem->rtmutex /* might block suspended (rt mutex) */
9:   clear READER BIAS
10:  if rw_sem->readers != 0 then
11:    suspend waiting for the last reader
12:  end if
13:  while 1 do
14:    take sem->rtmutex->wait_lock /* might block busy (spinlock) */
15:    if sem->readers == 0 then
16:      set WRITER BIAS
17:      release rw_sem->rtmutex->wait_lock
18:      return /* enter in the critical section */
19:    end if
20:    release rw_sem->rtmutex->wait_lock
21:    suspend waiting for the last reader
22:  end while
23:  return
24: end function

```

Fig. 4: Write-side operations

in Linux state-of-the-art protocols for (possibly heterogeneous) nested locks and developing novel analysis techniques. To the best of our knowledge, only few works on shared-memory multiprocessor synchronization targeted nested critical sections. Two notable examples are the work by Biondi et al. [5], in which a graph abstraction is introduced to derive a fine-grained analysis (i.e., not based on asymptotic bounds) for FIFO *non-preemptive* spin locks, and the one by Ward and Anderson [6], in which the *real-time nested locking protocol* (RNLP) is proposed, with the related *asymptotic* analysis. Later, Nemitz et al. [7] proposed an optimization for the average-case of RNLP. However, to the best of our knowledge, only the extension of RNLP proposed in [8] is explicitly conceived to deal with heterogeneous nested critical sections. The protocol is presented with the related asymptotic analysis, and an experimental study aimed at assessing schedulability. Future research work could target the issues in implementing the extended RNLP [8] in Linux. Also, it is worth considering the possibility of extending the graph abstraction proposed by Biondi et al. [5] to allow fine-grained analysis for nested heterogeneous locks.

The second open problem concerns the design of specialized analysis techniques accounting for specific implementations of complex types of locks (e.g., the aforementioned read/write lock in Linux). Considering the problem previously presented for the DOWN_WRITE function, an implementation-aware analysis would account for the contention on the heterogeneous nested critical section, considering it when a blocking-bound for the reader/writer semaphore is derived. The analyses for reader/writer semaphores that have already been proposed (e.g., the protocol proposed by Brandenburg and Anderson [9], or R/W RNLP [10], a variant of RNLP conceived to deal with nested, spin-based, read/write locks) could be integrated with implementation-specific aspects. The availability of blocking-bounds conceived considering the specific implementation adopted in the Linux kernel may help it to be more suitable for real-time contexts. Finally, a third open research area consists in finding more efficient locking protocols (with the related implementation), accounting for both general purpose benchmark performance (i.e., average-case behavior, needed by the GPOS nature of Linux) and predictability.

REFERENCES

- [1] J. Lelli, C. Scordino, L. Abeni, and D. Faggioli, "Deadline scheduling in the Linux kernel," *Software: Practice and Experience*, vol. 46, no. 6, pp. 821–839, 2016.
- [2] D. B. de Oliveira and R. S. de Oliveira, "Timing analysis of the PREEMPT RT linux kernel," *Softw., Pract. Exper.*, vol. 46, no. 6, pp. 789–819, 2016.
- [3] L. Sha, R. Rajkumar, and J. P. Lehoczky, "Priority inheritance protocols: An approach to real-time synchronization," *IEEE Transactions on Computers*, vol. 39, no. 9, Sep. 1990.
- [4] Red Hat, Inc., "Red Hat Enterprise Linux for Real Time," Available at: <https://www.redhat.com/it/resources/red-hat-enterprise-linux-real-time> [last accessed 28 March 2017].
- [5] A. Biondi, A. Weider, and B. Brandenburg, "A blocking bound for nested fifo spin locks," in *Real-Time Systems Symposium (RTSS)*, 2016, pp. 291–302.
- [6] B. C. Ward and J. H. Anderson, "Supporting nested locking in multiprocessor real-time systems," in *Real-Time Systems (ECRTS), 2012 24th Euromicro Conference on*, 2012, pp. 223–232.
- [7] C. E. Nemitz, T. Amert, and J. H. Anderson, "Real-time multiprocessor locks with nesting: Optimizing the common case," in *Proceedings of the 25th International Conference on Real-Time and Network Systems (RTNS 2017)*, 2017.
- [8] B. C. Ward and J. H. Anderson, "Fine-grained multiprocessor real-time locking with improved blocking," in *Proceedings of the 21st International Conference on Real-Time Networks and Systems*, ser. RTNS '13, 2013.
- [9] B. B. Brandenburg and J. H. Anderson, "Reader-writer synchronization for shared-memory multiprocessor real-time systems," in *2009 21st Euromicro Conference on Real-Time Systems*, July 2009, pp. 184–193.
- [10] B. C. Ward and J. H. Anderson, "Multi-resource real-time reader/writer locks for multiprocessors," in *2014 IEEE 28th International Parallel and Distributed Processing Symposium*, May 2014, pp. 177–186.