

Performance, Scalability, and Real-Time Response From the Linux Kernel

# Performance and Scalability Technologies in the Linux Kernel

#### Paul E. McKenney IBM Distinguished Engineer & CTO Linux Linux Technology Center



ACACES July 14, 2009

Copyright © 2009 IBM

#### **Course Objectives and Goals**

ß

- Introduction to Performance, Scalability, and Real-Time Issues on Modern Multicore Hardware: Is Parallel Programming Hard, And If So, Why?
- Performance and Scalability Technologies in the Linux Kernel
- Creating Performant and Scalable Linux Applications
- Real-Time Technologies in the Linux Kernel
   Creating Real-Time Linux Applications

#### Overview

8

# Programming Environments in Linux Kernel Synchronization Primitives

- Per-CPU Variables
- The Existence Problem

## Solutions to the Existence Problem

#### Overview

8

Programming Environments in Linux Kernel
Synchronization Primitives
Per-CPU Variables
The Existence Problem
Solutions to the Existence Problem
RCU API
Why Free and Open-Source Software?

## **Programming Environments in Linux Kernel**

 $\bigcirc$ 

#### **!PREEMPT Environments in Linux Kernel**



Ö

#### **PREEMPT** Environments in Linux Kernel



Ø

#### IBM

### PREEMPT\_RT Environments in Linux Kernel

8



ACACES 2009

## **Synchronization Primitives**

#### **Synchronization Primitives (Partial)**



#### **Synchronization Primitives (Partial)**



## Sequence Lock Example (Reader)

}

## Sequence Lock Example (Writer)

```
static inline void warp_clock(void)
{
    write_seqlock_irq(&xtime_lock);
    wall_to_monotonic.tv_sec -= sys_tz.tz_minuteswest * 60;
    xtime.tv_sec += sys_tz.tz_minuteswest * 60;
    update_xtime_cache(0);
    write_sequnlock_irq(&xtime_lock);
    clock_was_set();
}
```

Ø

#### **Sleep/Wakeup Primitives (Partial)**

- wait\_event(), wait\_event\_interruptible(), wait\_event\_timeout(), wait\_event\_killable(), wait\_event\_interruptible\_timeout()
  - \* wake\_up(), wake\_up\_process(), ...
- wait\_on\_bit()
  - \* wake\_up\_bit()

ACACES 2009

- These are favored over traditional sleep\_on() APIs for software-engineering reasons
  - \* "Unsafe at any speed" primitives get fixed
  - Sometimes repeatedly...

## **Synchronization Primitives**



#### But do this first!!!! Job #1 is *not* selecting primitives!

#### **Per-CPU Variables**

#### **Per-CPU Variables**

## DEFINE\_PER\_CPU(type, name) DECLARE\_PER\_CPU(type, name)

- per\_cpu(name, cpu)
- get\_cpu\_var(name)
- raw\_get\_cpu\_var(name)
- for\_each\_online\_cpu(var)
  - But careful!!! CPUs can come and go...
  - \* get\_online\_cpus() and put\_online\_cpus()

#### **The Existence Problem**

## **The Existence Problem**

The need for full partitioning suggests partitioning synchronization primitives, too!



http://www.usenix.org/events/osdi99/full\_papers/gamsa/gamsa.pdf

#### Access

 $\bigcirc$ 

- p = header;
- spin\_lock(&p->lock);
- do\_something\_with(p);
- spin\_unlock(&p->lock);

#### Deletion

- p = header;
- spin\_lock(&p->lock);
- if (header == p)

header = NULL;

#### else

```
p = NULL;
spin_unlock(&p->lock);
if (p != NULL)
    kfree(p);
```

If so, why? If not, why not, and what would be a solution?

#### Access

 $\bigcirc$ 

- Deletion
- p = header; if (p != NULL) { spin\_lock(&p->lock); do\_something\_with(p); spin\_unlock(&p->lock); }

```
q = p = header;
if (p != NULL) {
    spin_lock(&p->lock);
    if (header == p)
        header = NULL;
    else
        q = NULL;
    spin_unlock(&p->lock);
    if (q != NULL)
        kfree(q);
}
```

Never forget the NULL-pointer checks...

#### **The Existence Problem Extended**



 $\bigcirc$ 

#### Single-threaded programs

Locking chain extending back to "header"

- Similarly, transaction covering back to "header"
- Global lock

- Hashed global array of locks
- Per-CPU global locks
- Garbage collector
- Type-safe memory

ACACES 2009

- Global reference count
- Others?

#### Single-threaded programs

- Can be the right thing to do, but lose scalability
- And existence can be a problem even in singlethreaded code via signals, events, and callbacks

#### Locking/TM chain extending back to "header"

- High locking/STM overhead for nesting, increased probability of hitting HTM transaction-size limitations
- Deadlock issues with locking
- Increased probability of encompassing nonidempotent operation in both HTM and STM

#### Global lock

ACACES 2009

Poor scalability and performance

#### Hashed global array of locks

- Poor performance due to lack of memory locality
- Deadlock issues, especially for large programs
  - Engineering solutions well-known but complex

#### Per-CPU global locks

- Requires possibly-awkward partitioning over CPUs
- Deadlock issues, especially for large programs
  - Engineering solutions well-known but complex
- Type-safe memory
  - Complex code to detect and handle reallocation

#### Garbage collector

ACACES 2009

 Great if your environment provides one, but overhead and reclamation time can rule out GC

#### **Global Reference Count**



#### Access

8

#### Deletion

```
rcu read lock();
p = rcu dereference(header);
q = p;
if (p != NULL) {
    spin lock(&p->lock);
    if (header == p)
        header = NULL;
    else
         q = NULL;
    spin unlock(&p->lock);
rcu read unlock();
if (q != NULL) {
    synchronize rcu();
    kfree(q);
```

Reference acquired under rcu\_read\_lock() guaranteed to exist until rcu\_read\_unlock()

#### **Global Reference Count**

8



#### Does this work? Why or why not?

#### **Global Reference Count Issues**

- Updater starvation!!! (Why?)
- Horrible read-side performance under heavy contention (40us on 64-CPU Power-5 system)
- Mediocre read-side performance under light contention (100ns on Power-5 system)
- Extremely fast updates: 40 nanoseconds
  - \* But only in absence of readers

#### **Global Reference Count Pair Data**



#### **Global Reference Count Pair Data**

- 1 atomic\_t rcu\_refcnt[2];
- 2 atomic\_t rcu\_idx;

- 3 DEFINE\_SPINLOCK(rcu\_gp\_lock);
- 4 DEFINE PER CPU(int, rcu nesting);
- 5 DEFINE\_PER\_CPU(int, rcu\_read\_idx);

#### **Global Reference Count Pair Reader Primitives**

```
1 static void rcu read lock(void)
                                                             Acquire reference
 2 {
 3
     int i;
 4
     int n;
 5
     n = get cpu var(rcu nesting);
 6
     if (n == 0) {
 7
       i = atomic read(&rcu idx);
 8
         get cpu var(rcu read idx) = i;
 9
       atomic inc(&rcu refcnt[i]);
10
11
12
       get cpu var(rcu nesting) = n + 1;
13
     smp mb();
14 }
15
                                                             Release reference
16 static void rcu read unlock(void) 🗲
17 {
18
     int i;
19
     int n;
20
21
     smp mb();
     n = get cpu var(rcu nesting);
22
     if (n == 1) {
23
        i = get cpu var(rcu read idx);
24
25
       atomic dec(&rcu refcnt[i]);
26
27
       get cpu var(rcu nesting) = n - 1;
28 }
```

#### **Global Reference Count Pair Updater Primitives**

```
Wait for references
 1 void synchronize rcu(void)
                                                                 to be released
 2
     int i;
 3
 4
 5
     smp mb();
     spin lock(&rcu gp lock);
 6
     i = atomic read(&rcu idx);
 7
     atomic set(&rcu idx, !i);
 8
     smp mb();
 9
10
     while (atomic read(&rcu refcnt[i]) != 0) {
11
       poll(NULL, 0, 10);
12
13
     spin unlock(&rcu gp lock);
     smp mb();
14
15 }
```

#### Does this work? Why or why not?

#### **Issues With Global Reference Count Pair**

#### Buggy!!!

- CPU 0 rcu\_read\_lock() line 8: i = rcu\_idx == 0
- CPU 1 invokes synchronize\_rcu()
  - Now rcu\_idx == 1
- CPU 0 rcu\_read\_lock() line 9: atomically increments rcu\_refcnt[0], enters read-side critical section, acquires a reference to some data element
- CPU 1 removes that same data element
- CPU 1 invokes synchronize\_rcu() again:
  - Line 7 fetches i = rcu\_idx == 1
  - Line 8 sets rcu\_idx = 0
  - Lines 10-11 wait for rcu\_refcnt[1] to go to zero
- CPU 1 kfree()s the data element
  - While CPU 0 is still using it!!!

#### **Issues With Global Reference Count Pair**

#### Buggy!!!

- CPU 0 rcu\_read\_lock() line 8: i = rcu\_idx == 0
- \* CPU 1 invokes synchronize\_rcu()
  - Now rcu\_idx == 1
- CPU 0 rcu\_read\_lock() line 9: atomically increments rcu\_refcnt[0], enters read-side critical section, acquires a reference to some data element
- CPU 1 removes that same data element
- CPU 1 invokes synchronize\_rcu() again:
  - Line 7 fetches i = rcu\_idx == 1
  - Line 8 sets rcu\_idx = 0
  - Lines 10-11 wait for rcu\_refcnt[1] to go to zero
- CPU 1 kfree()s the data element
  - While CPU 0 is still using it!!!

Everyone who has attempted an implementation has committed this bug!

#### **Global Reference Count Pair Updater Primitives**



#### Does this work? Why or why not?

#### **Global Reference Count Pair Issues**

Horrible read-side performance under heavy contention (40us on 64-CPU Power-5 system) Double counter flip and update-side lock slow: 200ns in isolation, 40us on 64-CPU Power 5)

\* And no more concurrent updates!!!

- Mediocre-to-poor read-side performance under light contention (150ns on Power-5 system)
- No updater starvation

One thumb down!!!

#### **Per-CPU Reference Count?**

ACACES 2009

- Partitioning for performance and scalability
   But threads (tasks) take references, not CPUs!
  - So, disable preemption while holding reference
    - We are executing in the kernel, after all!!!
  - As is currently the case when holding spinlocks
- And when holding "raw" spinlocks in PREEMPT\_RT
   But if we are going to disable preemption...
  - \* The fact that the task is running on a given CPU is the reference!!!
  - Each context switch then implicitly releases the outgoing task's reference ...
  - \* ... and acquires the incoming task's reference

#### "Running on CPU" as Reference

```
1 static void rcu read lock()
                                                      Acquire reference
 2 {
     preempt disable(); /* no-op for !PREEMPT */
 3
 4 }
 5
                                                      Release reference
 6 static void rcu read unlock()
 7 {
     preempt enable(); /* no-op for !PREEMPT */
 8
 9 }
10
                                                      Wait for references
11 void synchronize rcu() 🗲
                                                        to be released
12 {
13
     int cpu;
14
     for each online cpu(cpu)
15
       run on(cpu);
16
17 }
```

Does this work? Why or why not? Can rcu\_read\_lock() participate in a deadlock cycle?

#### "Running on CPU" as Reference Evaluation

- Excellent read-side scalability and performance: "free" is a very good price!!!
- Scheduling on each CPU in turn works well for small numbers of CPUs, but does not scale well
  - \* Actual Linux-kernel implementation uses batching to achieve excellent update-side scalability
  - \* But is considerably more complex
- No updater starvation

ß



#### "Running on CPU" as Reference: Schematic



ß

#### "Running on CPU" as Reference: Use Case

#### Combines waiting for readers and multiple versions:

- Writer removes element B from the list (list\_del\_rcu())
- Writer waits for all readers to finish (synchronize\_rcu())
- Writer can then free B (kfree())



IBM

#### "Running on CPU" as Reference: Use Case

```
struct foo head {
                                  foo head
                                                     foo (A)
                                                                     foo (B)
                                                                                      foo (C)
   struct list head list;
   spinlock t mutex;
};
struct foo {
   struct list head list;
   int key;
                                                  int delete(struct foo head *fhp, int k)
};
                                                     struct foo *p;
int search(struct foo head *fhp, int k)
                                                     struct list head *head = &fhp->list;
   struct foo *p;
                                                     spin lock(&fhp->mutex);
   struct list head *head = &fhp->list;
                                                     list for each entry(p, head, list) {
                                                        if (p \rightarrow key == k) {
   rcu read lock();
                                                            list del rcu(p);
   list for each entry rcu(p, head, list)
                                                            spin unlock(&fhp->mutex);
      if (p \rightarrow key == k) {
                                                            synchronize rcu();
         rcu read unlock();
                                                            kfree(p);
         return 1;
                                                            return 1;
   }
   rcu read unlock();
                                                     spin unlock(&fhp->mutex);
   return 0;
                                                     return 0;
```

#### "Running on CPU" as Reference: rwlock

```
struct foo head {
   struct list head list;
                                  foo head
                                                     foo (A)
                                                                     foo (B)
                                                                                      foo (C)
   rwlock t mutex;
};
struct foo {
   struct list head list;
   int key;
                                                  int delete(struct foo head *fhp, int k)
};
                                                     struct foo *p;
int search(struct foo head *fhp, int k)
                                                     struct list head *head = &fhp->list;
{
   struct foo *p;
                                                     write lock(&fhp->mutex);
   struct list head *head = &fhp->list;
                                                     list for each entry(p, head, list) {
                                                        if (p \rightarrow key == k) {
   read lock(&fhp->mutex);
                                                            list del(p);
   list for each entry(p, head, list) {
                                                            write unlock(&fhp->mutex);
      if (p \rightarrow key == k) {
                                                            /* synchronize rcu(); */
         read lock(&fhp->mutex);
                                                            kfree(p);
         return 1;
                                                            return 1;
      }
   read unlock(&fhp->mutex);
                                                     write unlock(&fhp->mutex);
   return 0;
                                                     return 0;
}
                                                  }
```

#### **RCU vs. Reader-Writer Locking Performance**



#### CONFIG\_PREEMPT kernel build

#### EM

## What is RCU???

#### RCU is a:

- reader-writer lock replacement
- restricted reference-counting mechanism
- bulk reference-counting mechanism
- \* poor-man's garbage collector
- \* way of providing existence guarantees
- \* way of waiting for things to finish

#### Use RCU in:

- read-mostly situations or
- for deterministic response from read-side primitives and from asynchronous update-side primitives



#### **RCU** as a Solution to the Existence Problem

- For read-mostly data structures, RCU provides the benefits of the data-parallel model
  - But without the need to actually partition or replicate the RCU-protected data structures
  - Readers access data without needing to exclude each others or updates
    - Extremely lightweight read-side primitives
- And RCU provides additional read-side performance and scalability benefits
  - With a few limitations and restrictions....

ACACES 2009

#### **RCU for Read-Mostly Data Structures**

Ø

Almost...



RCU data-parallel approach: first partition resources, then partition work, and only then worry about parallel access control, and only for updates.

ACACES 2009

#### **RCU Usage in the Linux Kernel**

Ø



### **RCU Area of Applicability**

Read-Mostly, Stale & Inconsistent Data OK (RCU Works Great!!!)

Read-Mostly, Need Consistent Data (RCU Works OK)

Read-Write, Need Consistent Data (RCU *Might* Be OK...)

Update-Mostly, Need Consistent Data (RCU is *Really* Unlikely to be the Right Tool For The Job)



8

# Programming Environments in Linux Kernel Synchronization Primitives Per CBU Veriables

- Per-CPU Variables
- The Existence Problem

## Solutions to the Existence Problem

#### **Legal Statement**

ACACES 2009

 $\mathbb{A}$ 

- This work represents the view of the author and does not necessarily represent the view of IBM.
- IBM and IBM (logo) are trademarks or registered trademarks of International Business Machines Corporation in the United States and/or other countries.
- Linux is a registered trademark of Linus Torvalds.
- Other company, product, and service names may be trademarks or service marks of others.
- This material is based upon work supported by the National Science Foundation under Grant No. CNS-0719851.
  - Joint work with Manish Gupta, Maged Michael, Phil Howard, Joshua Triplett, and Jonathan Walpole

#### **Questions?**

To probe further:

- Linux Device Drivers, 3<sup>rd</sup> edition, J. Corbet, A. Rubini, G. Kroah-Hartman
- Linux Kernel Development, 2<sup>nd</sup> edition, Robert Love
- Linux Weekly News: lwn.net (Google for "whatever site:lwn.net")
- Linux Kernel source (http://kernel.org/pub/linux/kernel/v2.6/linux-2.6.30.tar.bz2)
  - "Documentation" directory
- http://lwn.net/Articles/262464/ (What is RCU, Fundamentally?)
- http://lwn.net/Articles/263130/ (What is RCU's Usage?)
- http://lwn.net/Articles/264090/ (What is RCU's API?)
- http://www.rdrop.com/users/paulmck/RCU/lockperf.2004.01.17a.pdf
  - Iinux.conf.au paper comparing RCU vs. locking performance
- http://www.rdrop.com/users/paulmck/RCU/RCUdissertation.2004.07.14e1.pdf
  - RCU motivation, implementations, usage patterns, performance (micro+sys)
- http://www.livejournal.com/users/james\_morris/2153.html
  - System-level performance for SELinux workload: >500x improvement
- http://www.rdrop.com/users/paulmck/RCU/hart\_ipdps06.pdf
  - Comparison of RCU and NBS (later appeared in JPDC)
- http://doi.acm.org/10.1145/1400097.1400099
  - History of RCU in Linux (Linux changed RCU more than vice versa)
- http://www.rdrop.com/users/paulmck/RCU/ (More RCU information)

#### **RCU API**



#### Why Free and Open-Source Software?

ð

#### Why Free and Open-Source Software?

#### The Parable of The Six Blind Penguins and the Elephant

ß

#### **Proprietary Programming: Requirements**



C

#### **Proprietary Programming: "Solution"**



#### But sooner or later...

Č

#### Example: DYNIX/ptx RCU Implementation

#### In late 1990s, knew everything there was to know about RCU:

- \* rcu\_read\_lock()
- \* rcu\_read\_unlock()
- \* call\_rcu()
- \* kfree()

8

- \* kmem\_cache\_free()
- \* kmem\_deferred\_free()

#### But DYNIX/ptx was only a database server...

#### IBM

#### The Entire Elephant Will Make Itself Known...



#### Which it did for RCU, but in Linux

![](_page_62_Picture_1.jpeg)

#### **FOSS Programming: Requirements**

![](_page_62_Picture_3.jpeg)

Ö

## Just Another Day on LKML...

Ó

![](_page_63_Picture_3.jpeg)

#### **But Sometimes Consensus is Achieved**

![](_page_64_Picture_3.jpeg)

#### Linux Community Taught Me Much About RCU

- Simplicity as a first-class requirement
- Support for 20+ CPU families
  - Must hide memory barriers from RCU users
- Real-time response: tens of microseconds
- Run on small-memory systems (2MB!!!)
- Heavy networking workloads
- Denial-of-service attacks

ACACES 2009

- Wait for interrupt- and NMI-handler completion
- Unloadable kernel modules
- Blocking in RCU read-side critical sections
- Sophisticated RCU-based list/tree manipulation
- Deep sub-millisecond RCU grace periods

![](_page_66_Picture_1.jpeg)

![](_page_66_Picture_2.jpeg)

Ö

### This is RCU

- rcu\_read\_lock()
- rcu\_read\_unlock()
- •
- .
- -

- -
- -
- call\_rcu()
- •
- •

- •
- •

- •
- •
- •
- •
- •
- kfree()
- kmem\_cache\_free()

## This is RCU in Linux

rcu\_read\_lock()

- rcu\_read\_unlock()
- rcu\_read\_lock\_bh()
- rcu\_read\_unlock\_bh()
- preempt\_disable()
- preempt\_enable()
- srcu\_read\_lock()
- srcu\_read\_unlock()
- rcu\_dereference()
- list\_for\_each\_entry\_rcu()
- hlist\_for\_each\_entry\_rcu()
- synchronize\_rcu()
- synchronize\_net()
- call\_rcu()
- rcu\_barrier()
- call\_rcu\_bh()

- synchronize\_sched()
- synchronize\_srcu()
- rcu\_assign\_pointer()
- list\_add\_rcu()
- list\_add\_tail\_rcu()
- list\_del\_rcu()
- list\_replace\_rcu()
- hlist\_del\_rcu()
- hlist\_add\_after\_rcu()
- hlist\_add\_before\_rcu()
- hlist\_add\_head\_rcu()
- hlist\_replace\_rcu()
- list\_splice\_init\_rcu()
- kfree()
- kmem\_cache\_free()

## This is RCU in Linux

rcu\_read\_lock()

8

- rcu\_read\_unlock()
- rcu\_read\_lock\_bh()
- rcu\_read\_unlock\_bh()
- preempt\_disable()
- preempt\_enable()
- srcu\_read\_lock()
- srcu\_read\_unlock()
- rcu\_dereference()
- list\_for\_each\_entry\_rcu()
- hlist\_for\_each\_entry\_rcu()

ACACES 2009

- synchronize\_rcu()
- synchronize\_net()
- call\_rcu()
- rcu\_barrier()
- call\_rcu\_bh()

- synchronize\_sched()
- synchronize\_srcu()
- rcu\_assign\_pointer()
- list\_add\_rcu()
- list\_add\_tail\_rcu()
- list\_del\_rcu()
- list\_replace\_rcu()
- hlist\_del\_rcu()
- hlist\_add\_after\_rcu()
- hlist\_add\_before\_rcu()
- hlist\_add\_head\_rcu()
- hlist\_replace\_rcu()
- list\_splice\_init\_rcu()
- kfree()
- kmem\_cache\_free()

#### **Any Questions?**