

---

---

# RCU in 2019

(from my perspective related to  
work I have been doing)

---

---

---

---

# Credits

RCU is the great decades-long work of Paul Mckenney and others. I am relatively new on the scene (~ 1.5 years).

---

---

# Who am I ; and how I got started with RCU?

- Worked on Linux for a decade or so.
- Heard only 5 people understand RCU. Opportunity !!!!
- Tired of reading RCU traces and logs that made no sense.
- Was running into RCU concepts but didn't know their meaning.

Time to put all mysteries to and end...

- Turned out I actually love it and it is like puzzle solving.

# Who am I ; and how I got started with RCU?

- Helping community / company with RCU issues, concepts, improvements, reviewing.
- Started coming up with RCU questions (~1.5 years ago) about why things could not be done a certain way to which PaulMck quickly reminded me that they are “correct” but... RCU is complex since it has to “survive” in the real-world.

# Who am I ; and how I got started with RCU?

PaulMck says...

“Here is your nice elegant  
little algorithm”



# Who am I ; and how I got started with RCU?

PaulMck says...

“Here is your nice elegant little algorithm equipped to survive in the Linux kernel”



# Agenda

- Introduction
- RCU Flavor consolidation
  - Performance
  - Scheduler Deadlock fixes
- List RCU API improvements (lockdep)

Extra material covered if time permitting:

- kfree\_rcu() batching
- Dynticks and NOHZ\_FULL fixes

# Intro: Typical RCU workflow

- Say you have some data that you have to share between a reader/writer section.

```
struct shared_data {  
    int a;  
    long b;  
};
```

```
int reader(struct shared_data *sd) {  
    if (sd->a)  
        return sd->b;  
    return 0;  
}
```

```
int writer(struct shared_data *sd) {  
    sd->b = 1;  
    sd->a = 2;  
}
```



# Intro: Typical RCU workflow

- One way is to use a reader-writer lock.

```
int reader(struct shared_data *sd) {  
    read_lock(&sd->rwlock);  
    if (sd->a)  
        ret = sd->b;  
    read_unlock(&sd->rwlock);  
    return ret;  
}
```

```
void writer(struct shared_data *sd) {  
    write_lock(&sd->rwlock);  
    sd->b = 1;  
    sd->a = 2;  
    write_unlock(&sd->rwlock);  
}
```

# Intro: Typical RCU workflow

- Or, use RCU:

```
int reader() {  
  
    struct shared_data sd =  
        rcu_dereference(global_sd);  
    rcu_read_unlock(); // READ  
    if (sd->a)  
        ret = sd->b;  
    rcu_read_unlock();  
  
    return ret;  
  
}
```

```
void writer() {  
    struct shared_data *sd, *old_sd;  
    spin_lock(&sd->lock);  
    old_sd = rcu_dereference(global_sd);  
    sd = kmalloc(sizeof(struct shared_data));  
    *sd = *old_sd; // COPY  
    sd->a = 2;  
    rcu_assign_pointer(global_sd, sd); // UPDATE  
    spin_unlock(&sd->lock);  
    synchronize_rcu();  
    kfree(old_sd);  
  
}
```

# Intro: Typical RCU workflow

- Writer makes a copy of data and modifies it.
- Readers execute in parallel on unmodified data.
- To know all prior readers are done, writer:
  - commit the modified copy (pointer)
  - Waits for a grace period.
  - After the wait, destroy old data.
  - Writer's wait marks the **start of a grace period**.
- Quiescent state (QS) concept - a state that the CPU passes through which signifies that the CPU is no longer in a read side critical section (tick, context switch, idle loop, use code, etc).
  - All CPUs must report passage through a QS.
  - Once all CPUs are accounted for, the **grace period can end**.

# Intro: What's RCU good for - summary

- Fastest Read-mostly synchronization primitive:
  - Readers are **almost free** regardless of number of concurrent ones. [Caveats]
  - Writes are costly but per-update cost is **amortized**.
    - 1000s or millions of updates can share a grace period, so cost of grace period cycle (updating shared state etc) is divided.
  - Read and Writer sections execute in **parallel**.

[Caveats about 'literally free']

- Slide 12: "Readers are almost free": Yes, but caveats do apply:
  - (1) CONFIG\_PREEMPT=n is free but =y has a small cost (especially unlock()).
  - (2) Possibility of rcu\_dereference() affecting compiler optimizations.
  - (3) Given that rcu\_read\_lock() and rcu\_read\_unlock() now imply barrier(), they can also affect compiler optimizations.

# Intro: What is RCU and what its RCU good for ?

- More use cases:
  - Wait for completion primitive (used in deadlock free NMI registration, BPF maps)
  - Reference counter replacements
  - Per-cpu reference counting
  - Optimized reader-writer locking (percpu\_rwsem)
    - Slow / Fast path switch (rcu-sync)
  - Go through RCUUsage paper for many many usage patterns.

# Toy #1 based on ClassicRCU (Docs: WhatIsRCU.txt)

Classic RCU (works only on PREEMPT=n kernels):

```
#define rcu_dereference(p) READ_ONCE(p);
#define rcu_assign_pointer(p, v) smp_store_release(&(p), (v));

void rcu_read_lock(void) { }

void rcu_read_unlock(void) { }

void synchronize_rcu(void)
{
    int cpu;
    for_each_possible_cpu(cpu)
        run_on(cpu);
}
```

---

---

## Only talking about TREE\_RCU today

TREE\_RCU is the most complex and widely used flavor of RCU.

*If you are claiming that I am worrying unnecessarily, you are probably right. But if I didn't worry unnecessarily, RCU wouldn't work at all!*

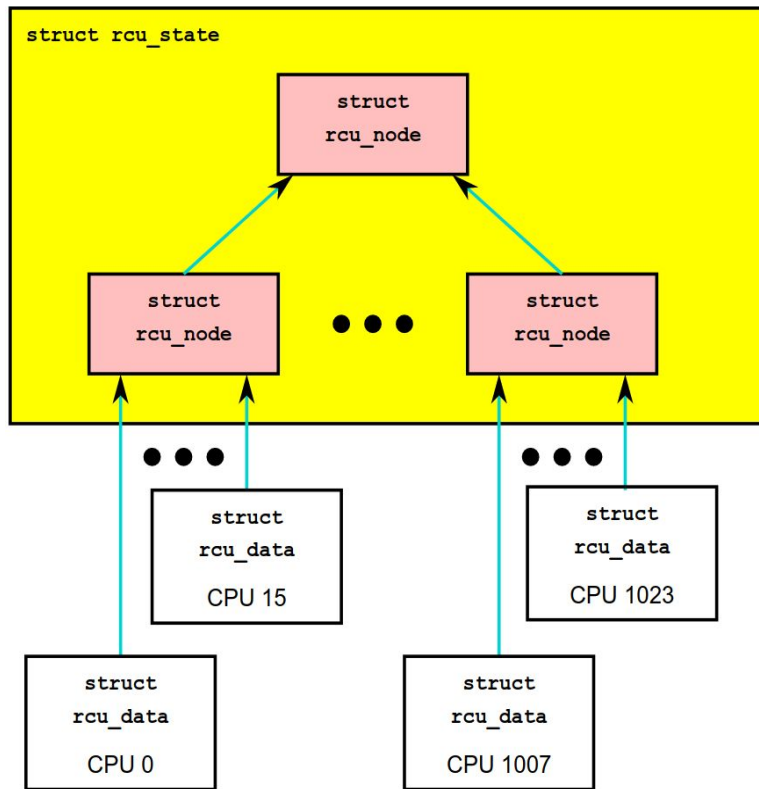
— Paul McKenney

---

---

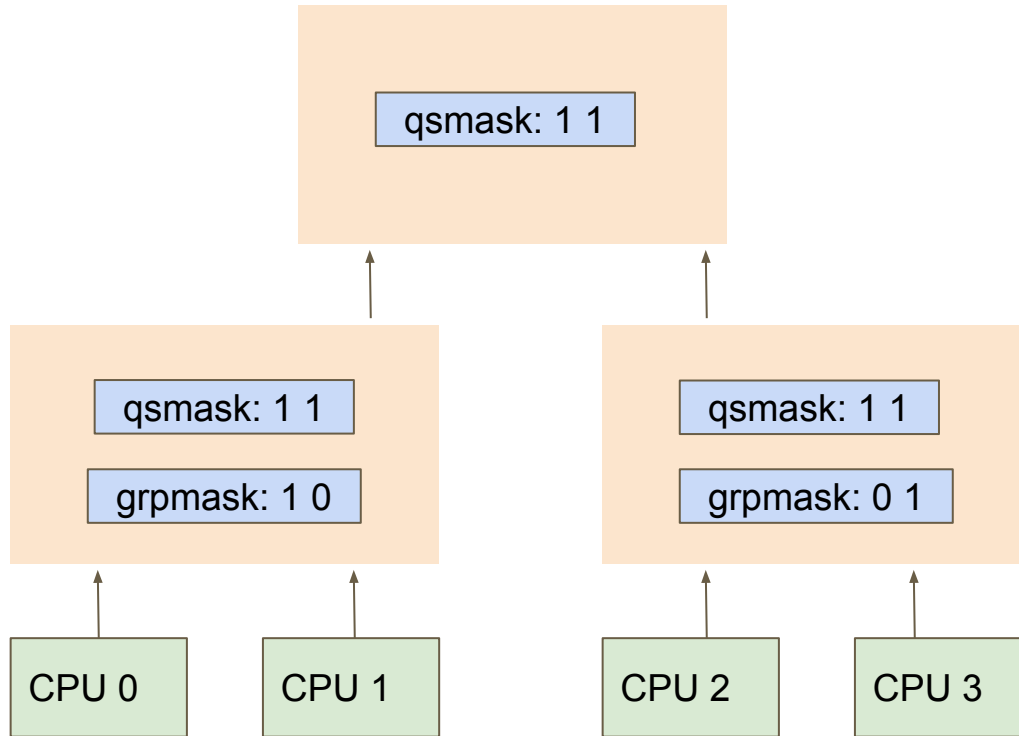
There's also TINY RCU but not discussing it.

# Intro: How TREE\_RCU works?

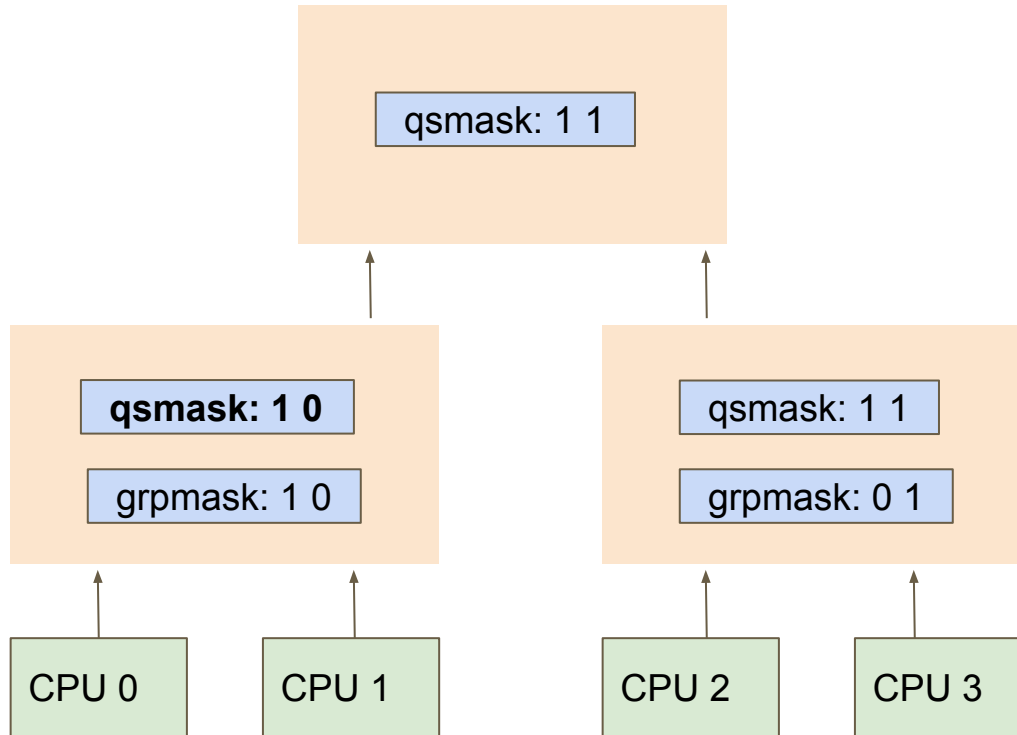




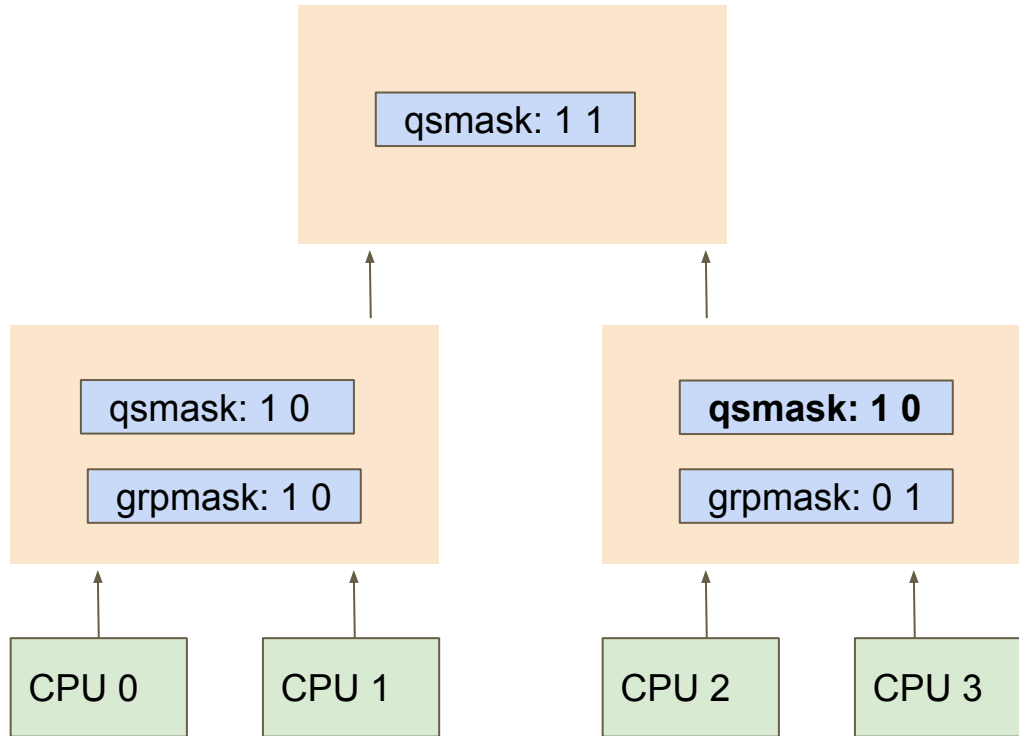
# TREE\_RCU example: Initial State of the tree



# TREE\_RCU example: CPU 1 reports QS

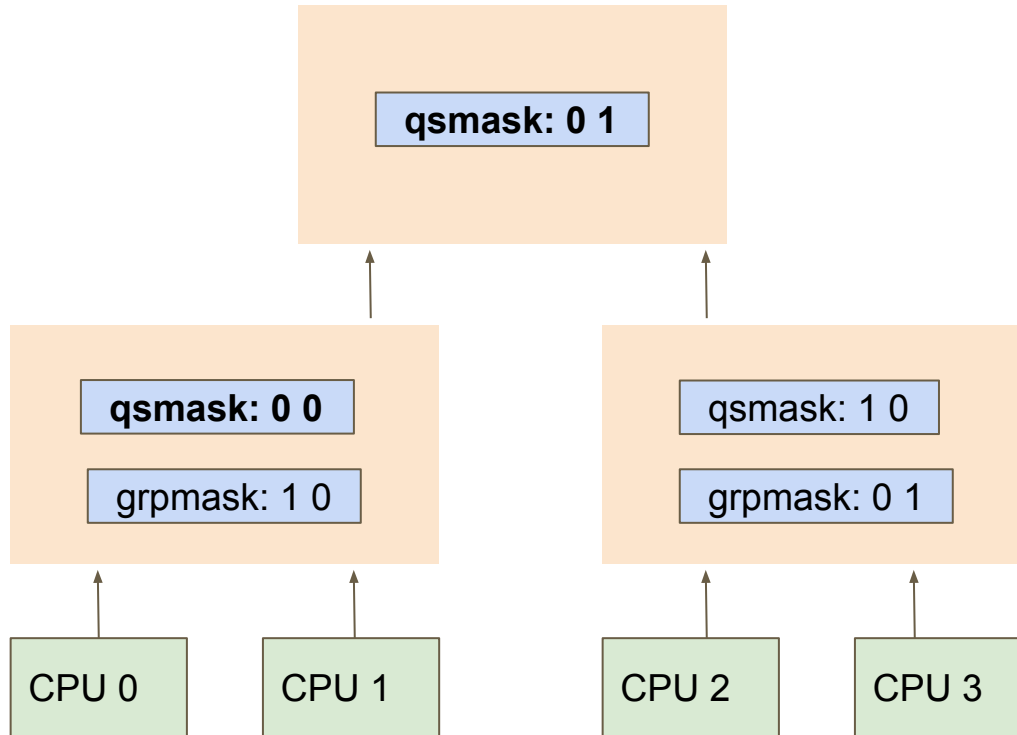


# TREE\_RCU example: CPU 3 reports QS



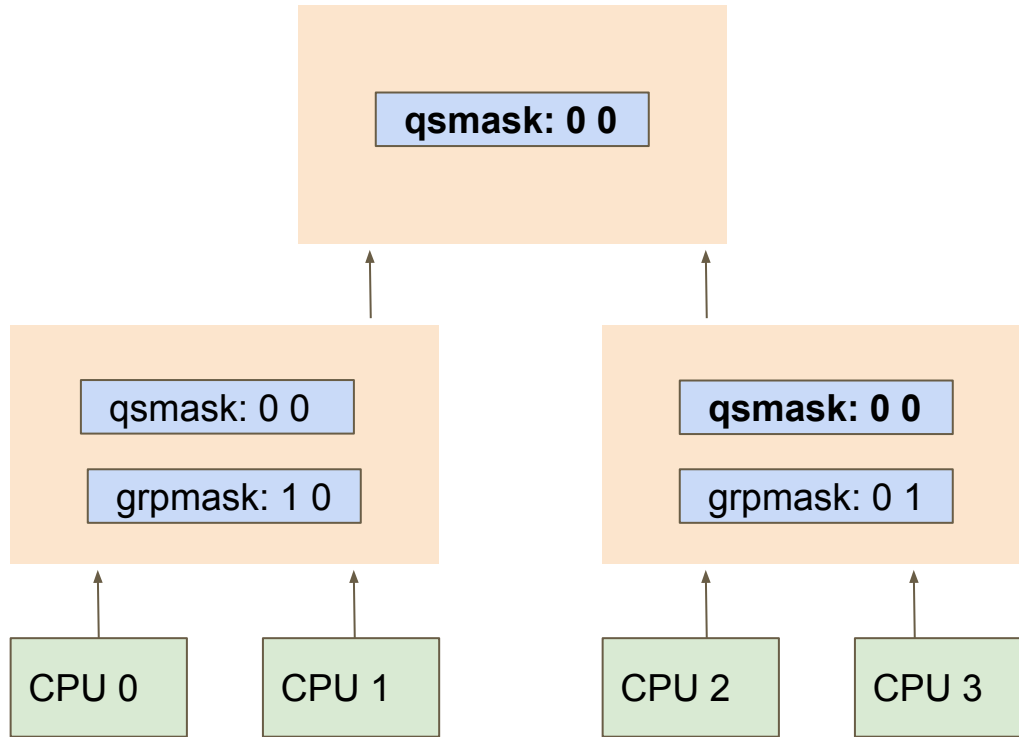
(Notice that the 2 QS updates have proceeded without any synchronization needed)

# TREE\_RCU example: CPU 0 reports QS



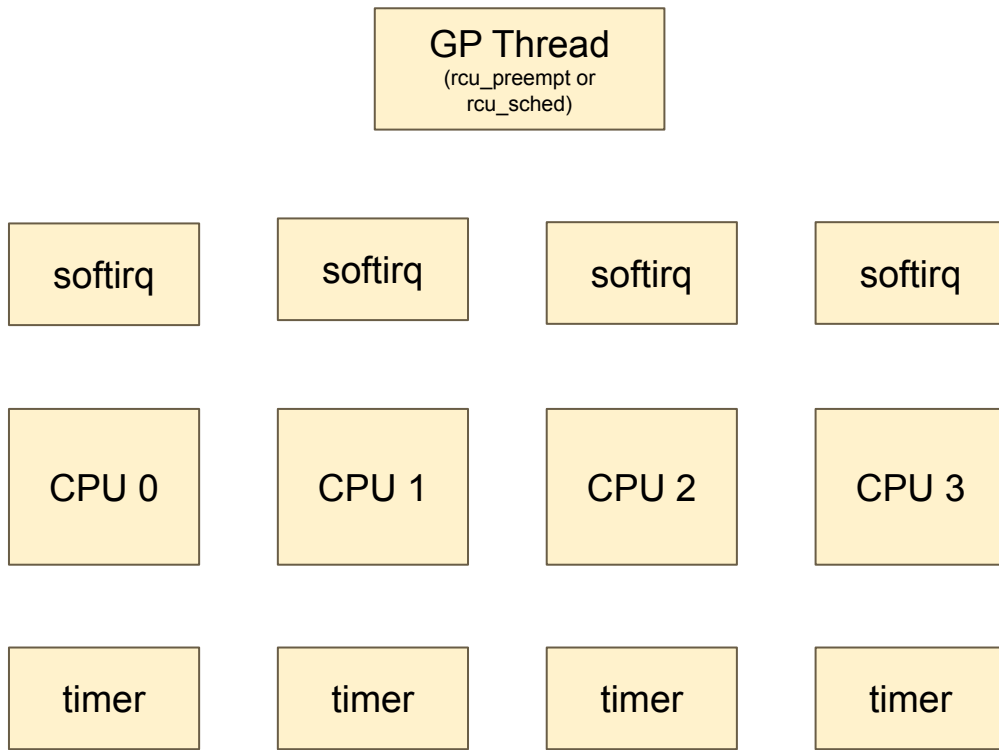
(Now there has been an update at the root node)

# TREE\_RCU example: CPU 2 reports QS

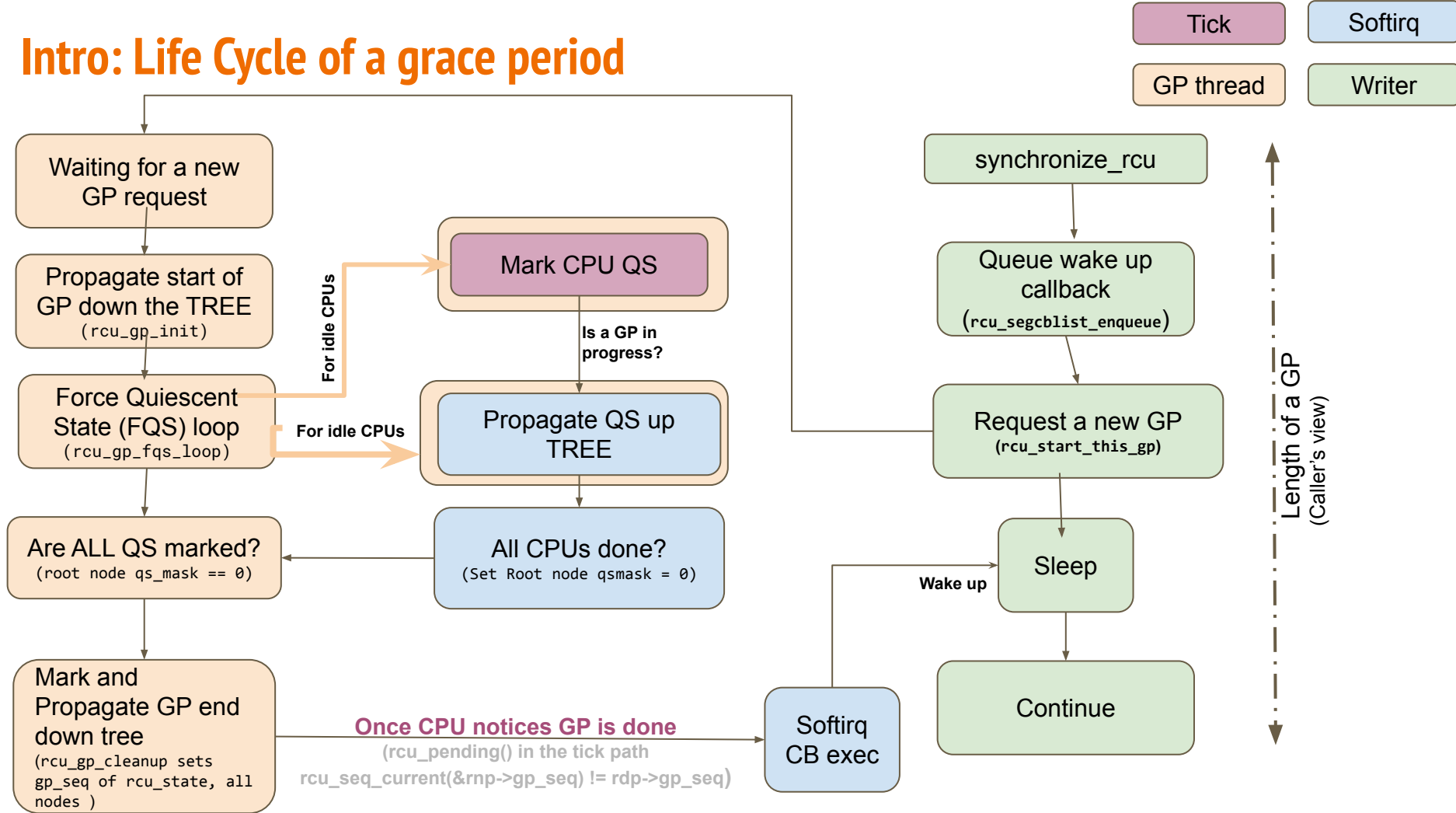


(Another update at the root node to mark qs has globally completed -- notice that only 2 global updates were needed instead of 4. On a system with 1000s of CPUs, this will be at most 64)

# Intro: Components of TREE RCU (normal grace period)



# Intro: Life Cycle of a grace period



# Intro: Where/When is Quiescent state reported

- Transition to / from CPU idle and user mode (a.k.a **dyntick-idle transitions**)
- **Context switch path** (If read-lock nesting is 0) -- NOTE: reports a CPU QS but task still blocks GP.
- In the Scheduler **tick**:
  - Task is not running in IRQ/preempt disable (no more sched/bh readers)
  - Read-lock nesting is 0 (no more preemptible rcu readers)
- In the FQS loop
  - Forcing quiescent states for **idle (nohz\_idle), usermode (nohz\_full), and offline** CPUs.
  - Forcing quiescent states for **CPUs that do dyntick-idle transitions.**
  - Including dyntick idle transitions faked with `rcu_momentary_dyntick_idle()` such as by `cond_resched()` on PREEMPT=n kernels.
  - Give soft hints to CPU's scheduler tick to resched task (`rdp.rcu_urgent_qs`)



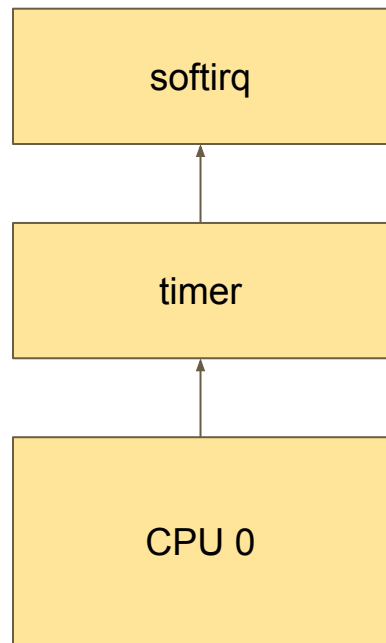
# Intro: Which paths do QS reports happen from?

- CPU local update is cheap:
  - Tick path.
  - Context Switch.
  - Help provided from `rcu_read_unlock()` if needed.
  - Reporting of a deferred QS reporting (when `rcu_read_unlock()` could not help).
- Global update of CPU QS -- Expensive, but REQUIRED for GP to end:
  - Softirq (once it sees CPU local QS from tick)
  - CPU going online (`rcu_cpu_starting`)
  - Grace period thread's FQS loop:
    - dyntick-idle QS report (Expensive, involves using atomic instructions)
      - dyntick-idle transitions (Kernel -> user ; Kernel -> idle)
      - `cond_resched()` in PREEMPT=n kernels does a forced dyntick-idle transition (HACK).

## Intro: What happens in softirq ?

### Per-CPU Work:

- **QS reporting for CPU and propagate up tree.**
- **Invoke any callbacks whose GP has completed.**



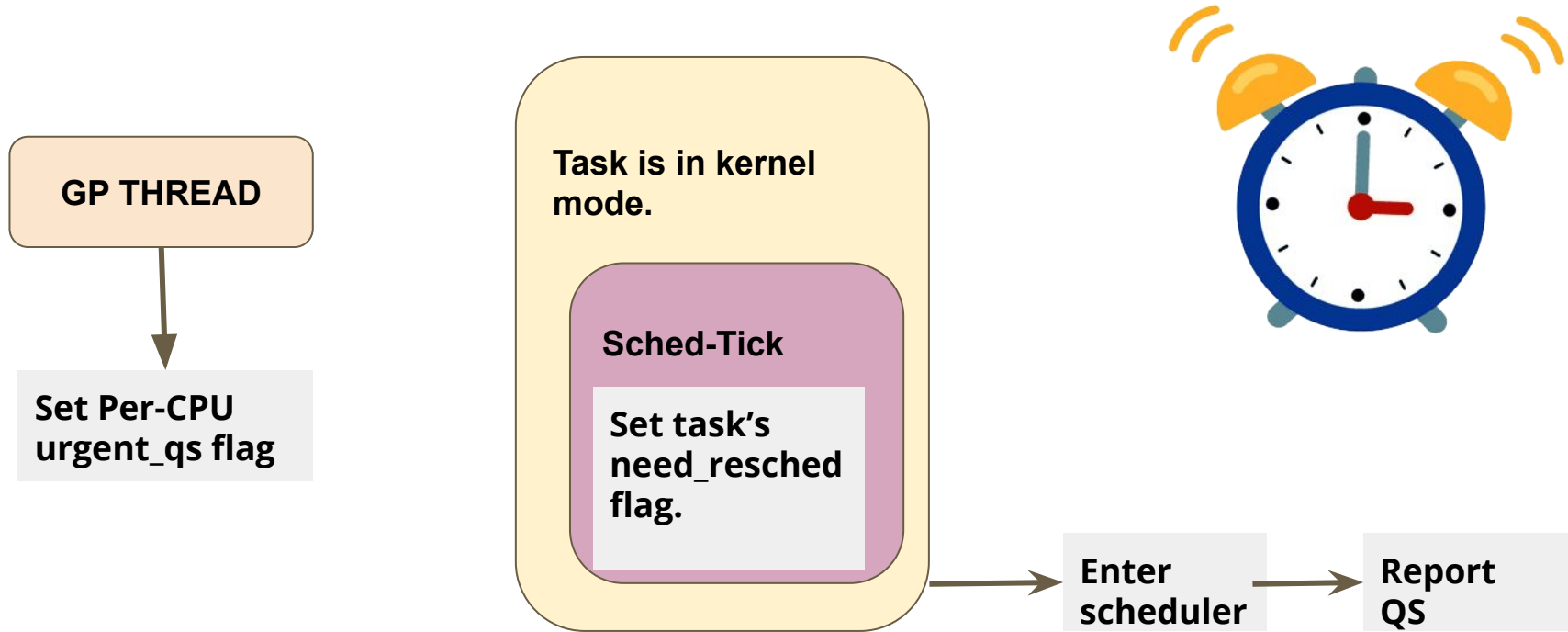
Caveat about callbacks queued on offline CPUs:

PaulMck says:

- > And yes, callbacks do migrate away from non-offloaded CPUs that go
- > offline. But that is not the common case outside of things like
- > rcutorture.

# Intro: Grace Period has started, what's RCU upto?

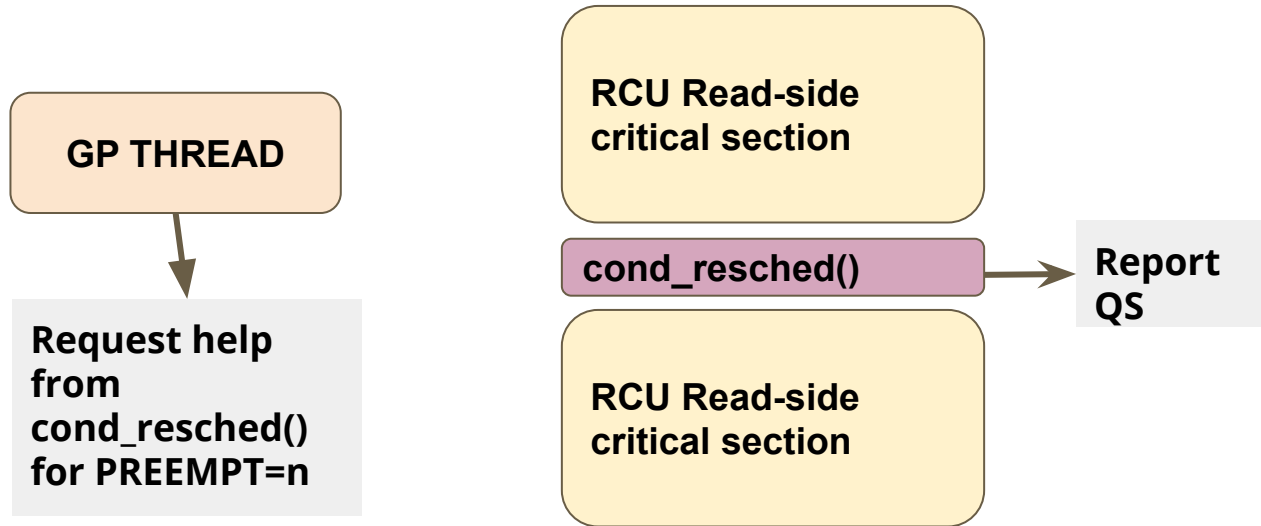
At around 100ms:



( Note: Scheduler entry can happen either in next TICK or next preempt\_enable() )

# Intro: Grace Period has started, what's RCU upto?

At around 200ms:

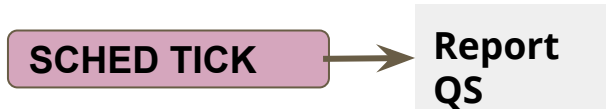
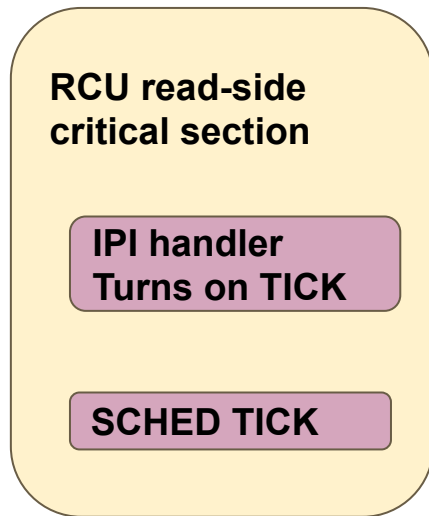
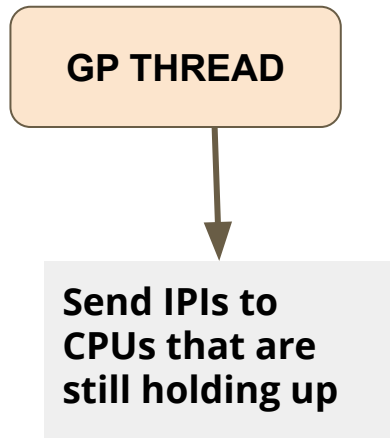


(by setting  
Per-cpu  
need\_heavy\_qs  
flag)



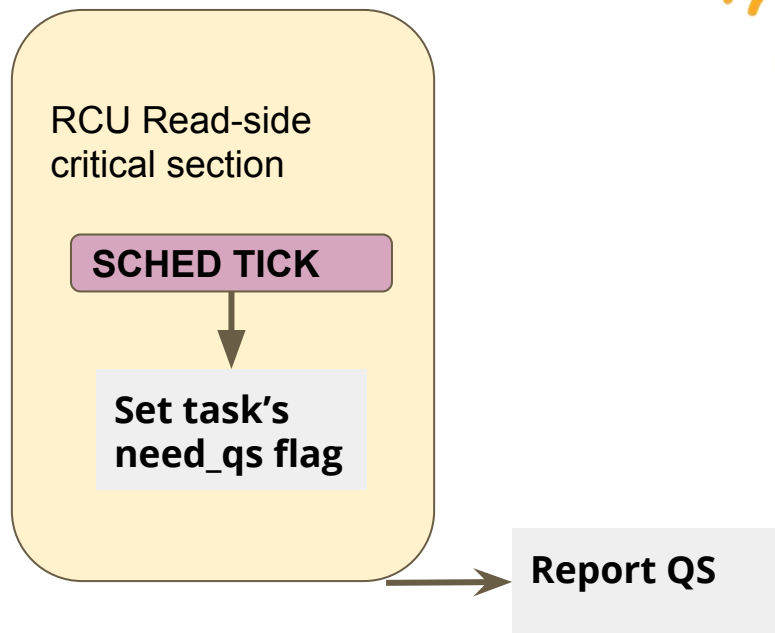
# Intro: Grace Period has started, what's RCU upto?

At around 300ms for NOHZ\_FULL CPUs:



# Intro: Grace Period has started, what's RCU upto?

At around 1 second of start of GP:



# Different RCU “flavors”

There are 3 main “flavors” of RCU in the kernel:

Entry into RCU read-side critical section:

1. “sched” flavor:
  - a. `rcu_read_lock_sched()`;
  - b. `preempt_disable()`;
  - c. `local_irq_disable()`;
  - d. IRQ entry.
  
2. “BH” flavor:
  - a. `rcu_read_lock_bh()`;
  - b. `local_bh_disable()`;
  - c. SoftIRQ entry.
  
3. “Preemptible” flavor only in CONFIG\_PREEMPT kernels
  - a. `rcu_read_lock()`;

# RCU Flavor Consolidation: Why? Reduce APIs

## Problem:

1. Too many APIs for synchronization. Confusion over which one to use!
  - a. For preempt flavor: `call_rcu()` and `synchronize_rcu()`.
  - b. For sched: `call_rcu_sched()` and `synchronize_rcu_sched()`.
  - c. For bh flavor: `call_rcu_bh()` and `call_rcu_bh()`.
2. Duplication of RCU state machine for each flavor ...

**Now after flavor consolidation:** Just `call_rcu()` and `synchronize_rcu()`.



# RCU Flavor Consolidation: Why? Changes to rcu\_state

3 rcu\_state structures before (v4.19)

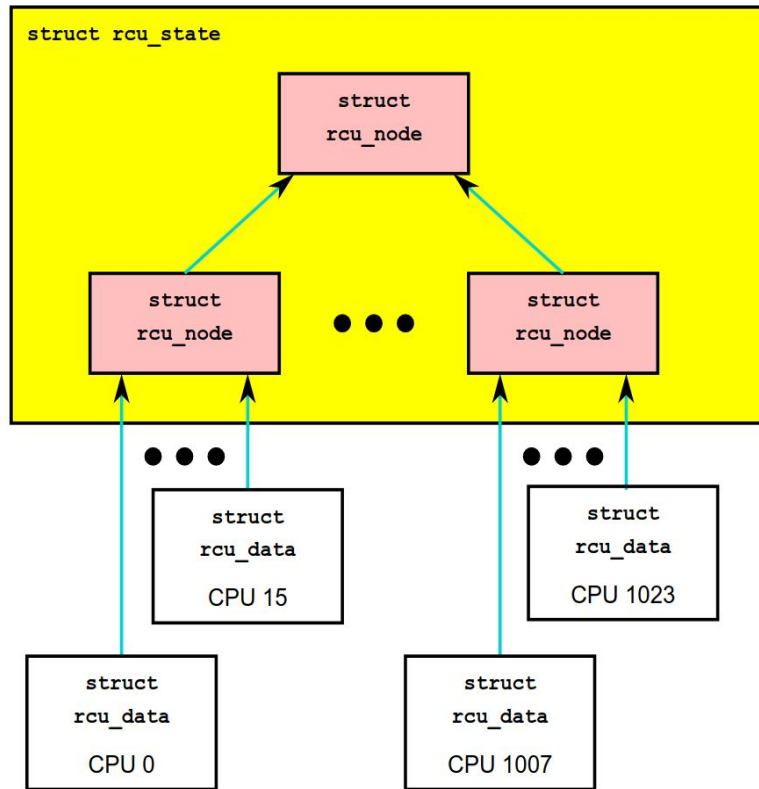
- rcu\_preempt
- rcu\_sched
- rcu\_bh

Now just 1 rcu\_state structure to handle all flavors:

- rcu\_preempt (Also handles sched and bh)

Number of GP threads also reduced from 3 to 1 since state machines have been merged now.

- Less resources!
- Less code!



# RCU Flavor Consolidation: Performance Changes

rcuperf test suite

- Starts N readers and N writers on N CPUs
- Each reader and writer affined to a CPU
- Readers just do `rcu_read_lock/rcu_read_unlock` in a loop
- Writers start and end grace periods by calling `synchronize_rcu` repeatedly.
- Writers measure time before/after calling `synchronize_rcu`.

**Locally modified test to do on reserved CPU continuously:**

`preempt_disable + busy_wait + preempt_enable` in a loop

Note: This is an additional reader-section variant now (sched-RCU) running in parallel to the existing `rcu_read_lock` readers (preemptible-RCU).

**What could be the expected Results?**

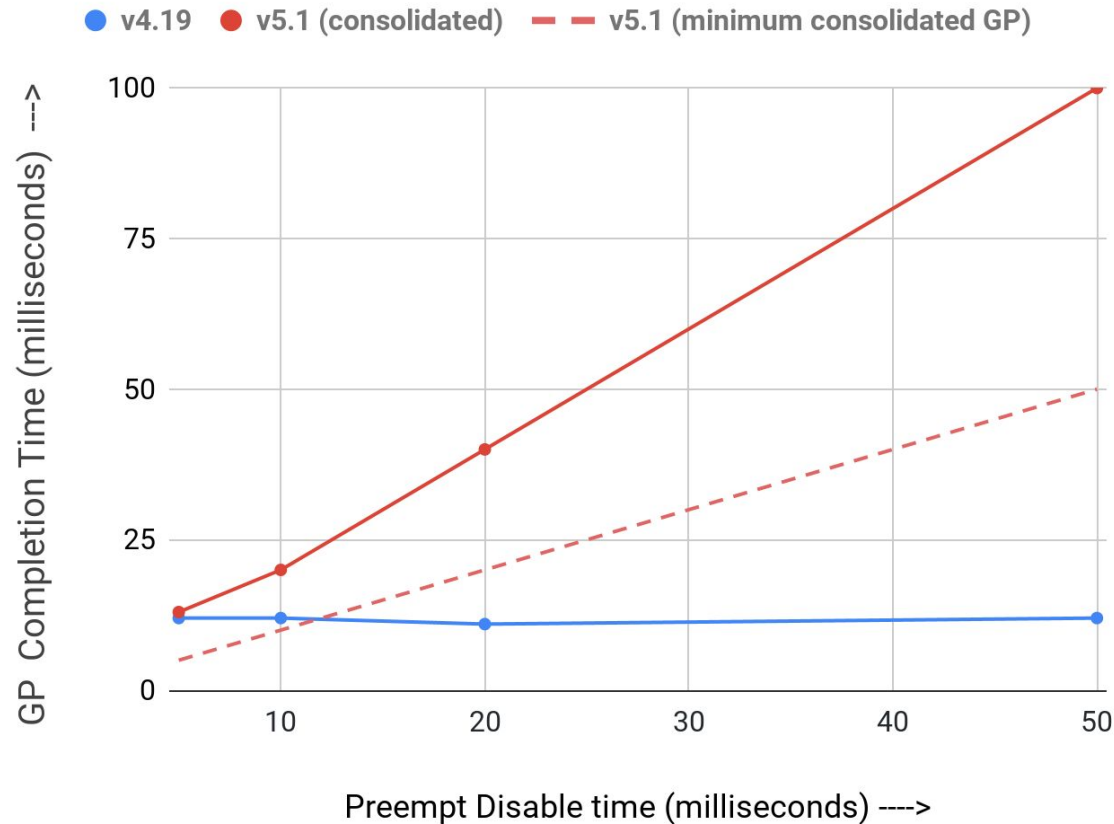
# RCU Flavor Consolidation

## Performance Changes

This is still within RCU specification!

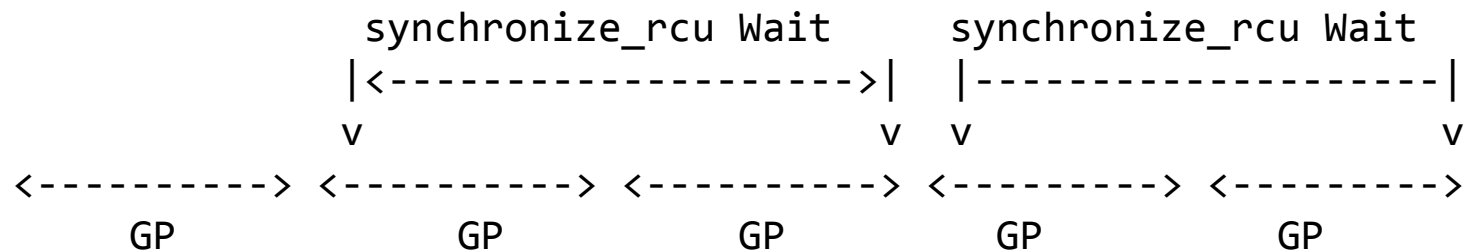
Also note that disabling preemption for so long is most not acceptable by most people anyway.

Comparison of v4.19 and v5.1 with rcuperf mods



# RCU Flavor Consolidation

Notice that `synchronize_rcu` time was 2x the `preempt_disable` time, that's cos:



GP = long preempt disable duration

# Consolidated RCU - The different cases to handle

Say RCU requested special help from the reader section unlock that is holding up a GP for too long....

```
preempt_disable();
rcu_read_lock();
do_some_long_activity(); // TICK sets per-task ->need_qs bit
rcu_read_unlock();      // Now need special processing from
                        // rcu_read_unlock_special()
preempt_enable();
```

## RCU-preempt reader nested in RCU-sched

## Consolidated RCU - The different cases to handle

### Before:

```
preempt_disable();
rcu_read_lock();
do_some_long_activity();
rcu_read_unlock()
    -> rcu_read_unlock_special(); // Report the QS
preempt_enable();
```

### Now:

```
preempt_disable();
rcu_read_lock();
do_some_long_activity();
rcu_read_unlock();
    -> rcu_read_unlock_special(); // Defer the QS and set
                                   // rcu_read_unlock_special.deferred_qs
                                   // bit & set TIF_NEED_RESCHED
preempt_enable(); // Report the QS
```

## RCU-preempt reader nested in RCU-bh

## Consolidated RCU - The different cases to handle

### Before:

```
local_bh_disable(); /* or softirq entry */
rcu_read_lock();
do_some_long_activity();
rcu_read_unlock()
    -> rcu_read_unlock_special(); // Report the QS
local_bh_enable(); /* softirq exit */
```

### Now:

```
local_bh_disable(); /* or softirq entry */
rcu_read_lock();
do_some_long_activity();
rcu_read_unlock();
    -> rcu_read_unlock_special(); // Defer the QS and set
                                   // rcu_read_unlock_special.deferred_qs
                                   // bit & set TIF_NEED_RESCHED
local_bh_enable(); /* softirq exit */ // Report the QS
```

# Consolidated RCU - The different cases to handle

**RCU-preempt reader nested in local\_irq\_disable (IRQoff)**  
(This is a special case where previous reader requested deferred special processing by setting ->deferred\_qs bit)

**Before:**

```
local_irq_disable();
rcu_read_lock();
rcu_read_unlock()
    -> rcu_read_unlock_special(); // Report the QS
local_irq_enable();
```

**Now:**

```
local_irq_disable();
rcu_read_lock();
rcu_read_unlock();
    -> rcu_read_unlock_special(); // Defer the QS and set
                                   // rcu_read_unlock_special.deferred_qs
                                   // bit & set TIF_NEED_RESCHED
local_irq_enable(); // CANNOT Report the QS, still deferred.
```



# Consolidated RCU - The different cases to handle

## RCU-preempt reader nested in RCU-sched (IRQoff)

(This is a special case where previous reader requested deferred special processing by setting ->deferred\_qs bit)

### Before:

```
/* hardirq entry */
rcu_read_lock();
rcu_read_unlock()
    -> rcu_read_unlock_special(); // Report the QS
/* hardirq exit */
```

### Now:

```
/* hardirq entry */
rcu_read_lock();
do_some_long_activity();
rcu_read_unlock();
    -> rcu_read_unlock_special(); // Defer the QS and set
                                   // rcu_read_unlock_special.deferred_qs
                                   // bit & set TIF_NEED_RESCHED
/* hardirq exit */    // Report the QS
```

## RCU-bh reader nested in RCU-preempt reader

# Consolidated RCU - The different cases to handle

### Before:

```
/* hardirq entry */
rcu_read_lock();
rcu_read_unlock()
    -> rcu_read_unlock_special(); // Report the QS
/* hardirq exit */
```

### Now:

```
/* hardirq entry */
rcu_read_lock();
do_some_long_activity();
rcu_read_unlock();
    -> rcu_read_unlock_special(); // Defer the QS and set
                                   // rcu_read_unlock_special.deferred_qs
                                   // bit & set TIF_NEED_RESCHED
/* hardirq exit */ // Report the QS
```

## RCU-bh reader nested in RCU-preempt or RCU-sched

### Before:

```
preempt_disable();
/* Interrupt arrives */
/* Raises softirq */
/* Interrupt exits */
__do_softirq();
    -> rcu_bh_qs();          /* Reports a BH QS */
preempt_enable();
```

### Now:

```
preempt_disable();
/* Interrupt arrives */
/* Raises softirq */
/* Interrupt exits */
__do_softirq();          /* Do nothing -- preemption still disabled */
preempt_enable();
```

## Consolidated RCU - The different cases to handle

# Consolidated RCU - The different cases to handle

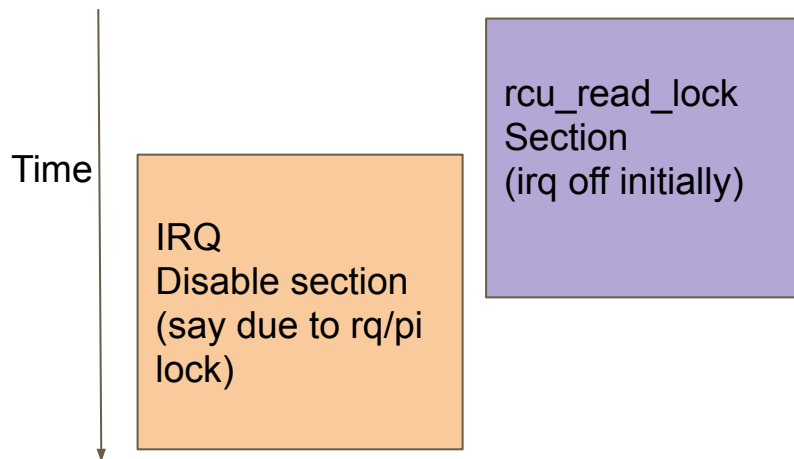
**Solution:** In case of denial of attack, ksoftirqd's loop with report QS.  
No reader sections expected there:

See commit: d28139c4e967 ("rcu: Apply RCU-bh Qses to RCU-sched and RCU-preempt when safe")

# Consolidated RCU - Fixing scheduler deadlocks...

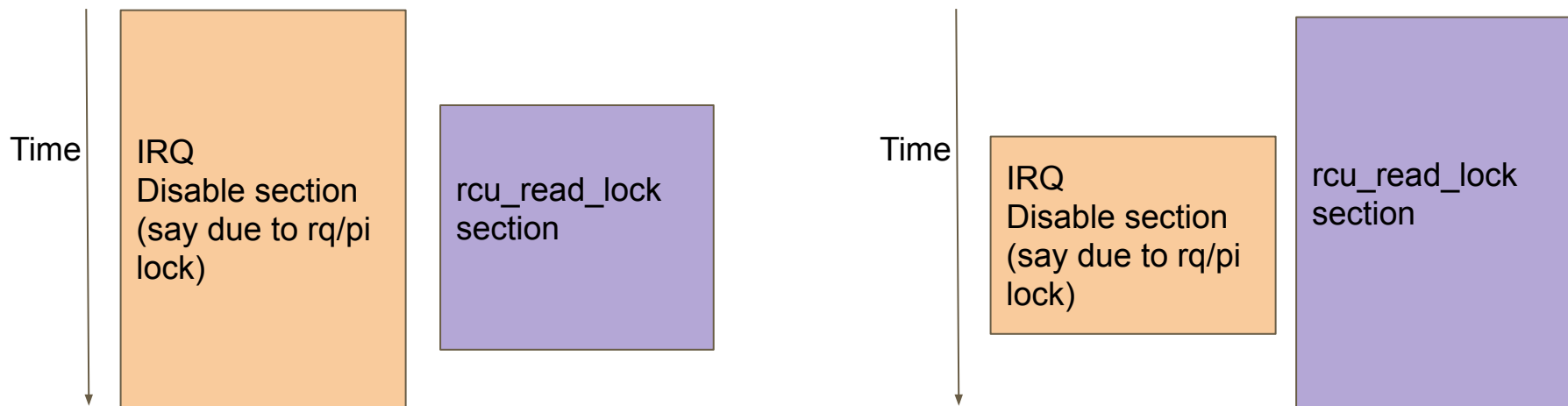
The forbidden scheduler rule... This is NOT allowed (<https://lwn.net/Articles/453002/>)

“Thou shall not hold RQ/PI locks across an rcu\_read\_unlock() if thou not holding it or disabling IRQ across both both the rcu\_read\_lock() + rcu\_read\_unlock()”



# Consolidated RCU - Fixing scheduler deadlocks...

The forbidden scheduler rule... This is allowed



# Consolidated RCU - Fixing scheduler deadlocks...

But we have a new problem... Consider case: future `rcu_read_unlock_special()` might be called due to a previous one being deferred.

```
previous_reader()
{
    rcu_read_lock();
    do_something(); /* Preemption happened here (so need help from rcu_read_unlock_special. */
    local_irq_disable(); /* Cannot be the scheduler as we discussed! */
    do_something_else();
    rcu_read_unlock(); // As IRQs are off, defer QS report but set deferred_qs bit in rcu_read_unlock_special
    do_some_other_thing();
    local_irq_enable();
}

current_reader() /* QS from previous_reader() is still deferred. */
{
    local_irq_disable(); /* Might be the scheduler. */
    do_whatever();
    rcu_read_lock();
    do_whatever_else();
    rcu_read_unlock(); /* Must still defer reporting QS once again but safely! */
    do_whatever_comes_to_mind();
    local_irq_enable();
}
```

# Consolidated RCU - Fixing scheduler deadlocks...

**Fixed in commit: 23634eb ("rcu: Check for wakeup-safe conditions in rcu\_read\_unlock\_special()")**

Solution: Intro `rcu_read_unlock_special.b.deferred_qs` bit. (Which is set in `previous_reader()` in previous example).

Raise `softirq` from `_special()` only when either of following are true:

- `in_irq()` (later changed to `in_interrupt`) - because `ksoftirqd` wake-up impossible.
- `deferred_qs` is set which happens in `previous_reader()` in previous example.

This makes the `softirq` raising not wake `ksoftirqd` thus avoiding a scheduler deadlock.

**Made detailed notes on scheduler deadlocks:**

<https://people.kernel.org/joelfernandes/making-sense-of-scheduler-deadlocks-in-rcu>

<https://lwn.net/Articles/453002/>



# CONFIG\_PROVE\_RCU does “built-in” deref check

```
rcu_read_lock();
```

```
// rcu_dereference checks if rcu_read_lock was called  
x = rcu_dereference(y);
```

```
rcu_read_unlock();
```

# However this can mislead it...

```
spin_lock(&my_lock);
```

```
// Perfectly legal but still splats!
```

```
x = rcu_dereference(y);
```

```
spin_unlock(&my_lock);
```

# So there's an API to support such dereference:

```
spin_lock(&my_lock);
```

```
x = rcu_dereference_protected(y, lockdep_is_held(&my_lock));
```

```
spin_unlock(&my_lock);
```

# List RCU usage - hardening

- listRCU APIs: `list_for_each_entry_rcu()` don't do any checking.
- Can hide subtle RCU related bugs such as list corruption.

# List RCU usage - Example of issue, consider...

```
/*
 * NOTE: acpi_map_lookup() must be called with rcu_read_lock() or spinlock protection.
 * But nothing checks for that, list corruption will go unnoticed...
 */

struct acpi_ioremap * acpi_map_lookup(acpi_physical_address phys, acpi_size size)
{
    struct acpi_ioremap *map;
    list_for_each_entry_rcu(map, &acpi_ioremaps, list, acpi_ioremap_lock_held())
        if (map->phys <= phys &&
            phys + size <= map->phys + map->size)
            return map;
}
```

# List RCU usage - hardening

- Option 1: Introduce new APIs that do checking for each RCU “flavor” :
  - `list_for_each_entry_rcu()`
  - `list_for_each_entry_rcu_bh()`
  - `list_for_each_entry_rcu_sched()`
- And for when locking protection is used instead or reader protection:
  - `list_for_each_entry_rcu_protected()`
  - `list_for_each_entry_rcu_bh_protected()`
  - `list_for_each_entry_rcu_sched_protected()`
- Drawback:
  - Proliferation of APIs...

# List RCU usage - hardening

- Option 2: Can do better since RCU backend is now consolidated.
  - So checking whether ANY flavor of RCU reader is held, is sufficient in `list_for_each_entry_rcu()`
    - `sched`
    - `bh`
    - preemptible RCU
  - And add an optional argument to `list_for_each_entry_rcu()` to pass a lockdep expression; so no need for a `list_for_each_entry_rcu_protected()`.
- One API to rule it all !!!!!

# List RCU usage - hardening

Example usage showing optional expression when lock held: `acpi_map_lookup()`:

```
#define acpi_ioremap_lock_held() lockdep_is_held(&acpi_ioremap_lock)

struct acpi_ioremap * acpi_map_lookup(acpi_physical_address phys, acpi_size size)
{
    struct acpi_ioremap *map;
    list_for_each_entry_rcu(map, &acpi_ioremaps, list, acpi_ioremap_lock_held())
        if (map->phys <= phys &&
            phys + size <= map->phys + map->size)
            return map;
}
```



# List RCU usage - hardening

Patches:

<https://lore.kernel.org/lkml/20190716184656.GK14271@linux.ibm.com/T/#t>

<https://lore.kernel.org/patchwork/project/lkml/list/?series=401865>

Next steps:

- Make CONFIG\_PROVE\_RCU\_LIST the default and get some more testing

# Future work

- More Torture testing on arm64 hardware
- Re-design dynticks counters to keep simple
- List RCU checking updates
- RCU scheduler deadlock checking
- Reducing grace periods due to `kfree_rcu()`.
- Make possible to not embed `rcu_head` in object
- More RCU testing, experiment with modeling etc.
- More systematic study of `__rcu` sparse checking.

# Thank you ...

- For questions, please email the list: [rcu@vger.kernel.org](mailto:rcu@vger.kernel.org)
- Follow us on Twitter:
  - paulmckrcu@ joel\_linux@ boqun\_feng@ srostedt@

# kfree\_rcu() improvements

Problem: Too many kfree\_rcu() calls can overwhelm RCU

- Concerns around vhost driver calling kfree\_rcu often:
  - <https://lore.kernel.org/lkml/20190721131725.GR14271@linux.ibm.com/>
  - <https://lore.kernel.org/lkml/20190723035725-mutt-send-email-mst@kernel.org/>
- Similar bug was fixed by avoiding kfree\_rcu() altogether in access()
  - <https://lore.kernel.org/lkml/20190729190816.523626281@linuxfoundation.org/>

Can we improve RCU to do better and avoid future such issues?

# Current kfree\_rcu() design (TODO: Add state diagram)

kfree\_rcu -> queue call back -> wait for GP -> run kfree from RCU core

- Flooding kfree\_rcu() can cause lot of RCU work and grace periods.
- All kfree() happens from softirq as RCU callbacks can execute from softirq
  - Unless nocb or nohz\_full is passed of course, still running more callbacks means other RCU work also has to wait.

# kfree\_rcu() re-design - first pass

kfree\_rcu batching:

- 2 per-CPU list (Call these A and B)
- Each kfree\_rcu() request queued on a per-CPU list (no RCU)
- Within KFREE\_MAX\_JIFFIES (using schedule\_delayed\_work())
  - Dequeue all requests on A, and Queue on B
  - Call queue\_rcu\_work() to run a worker to free all B-objects after GP
  - Meanwhile new requests can be queued on A

# kfree\_rcu() re-design - first pass

## Advantages:

- 5X reduction in number of grace periods in rcuperf testing:
  - GPs from **~9000 to ~1700** (per-cpu kfree\_rcu flood on 16 CPU system)
- All kfree() happens in workqueue context and systems schedules work on the different CPUs.

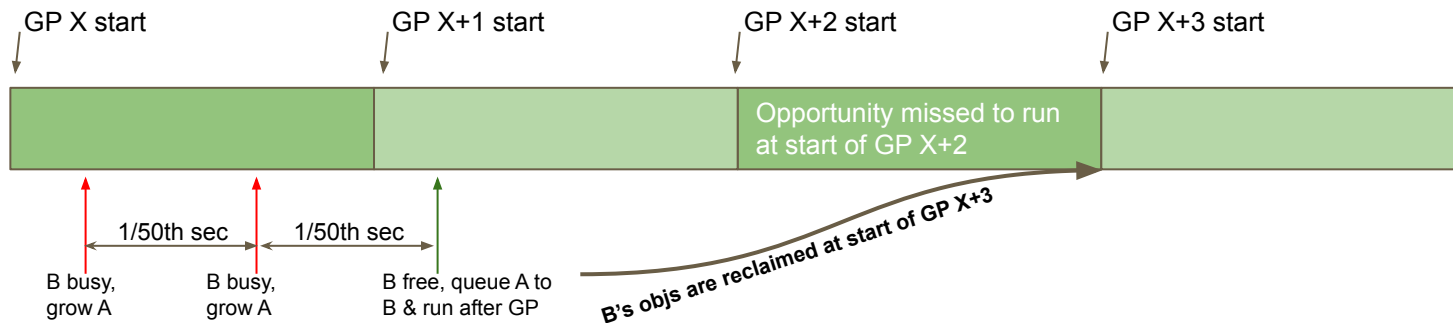
## Drawbacks:

- Since list B cannot be queued into until GP ends, list A keeps growing
- Causes high memory consumption (300-350) MB.

# kfree\_rcu() re-design - first pass

Drawbacks:

- Since list B cannot be queued into until GP, list A keeps growing
- Miss opportunities and has to wait for future GPs on busy system
- Causes higher memory consumption (300-350) MB.



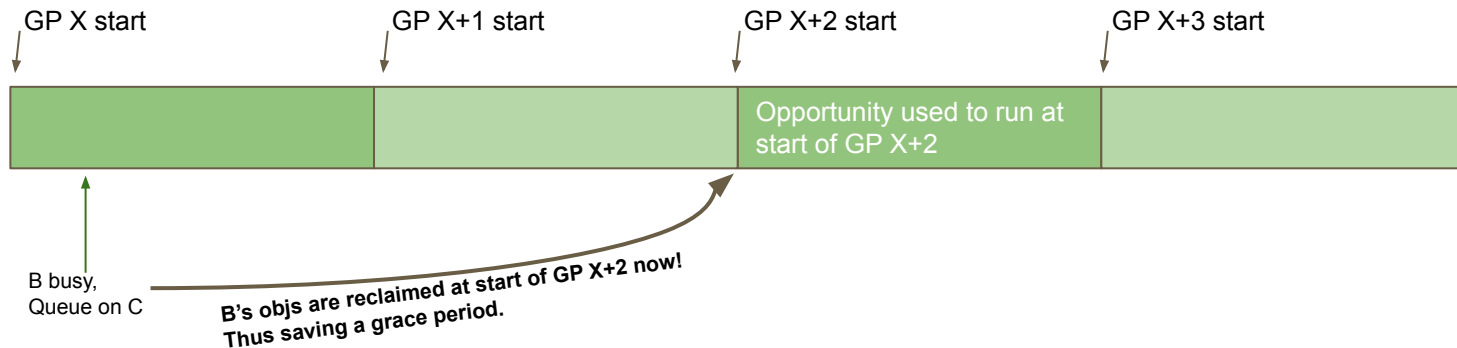
B is busy, so any attempts to Queue A to B have to be delayed.

So we keep retrying every 1/50th of a second; and A keeps growing...



# kfree\_rcu() re-design - second pass

- Introduce a new list C.
  - If B is busy just queue on C and schedule that after GP



- Brings down memory consumption by 100-150MB during kfree\_rcu flood
- Maximize bang for GP-bucks!

# dyntick RCU counters : intro (TODO: Condense these 9 slides into one slide)

- RCU needs to know what state CPUs are in
- If we have scheduler tick, no problem! Tick path tells us.

But...

- NOHZ\_FULL can turn off tick in both user mode and idle mode.
- NOHZ\_IDLE can turn off tick in idle mode

Implemented with a per-CPU atomic counter:

- `rcu_data :: atomic_t dynticks` : Even for RCU-idle, odd non-RCU-idle

# dyntick RCU counters : intro (TODO: Add state-diagrams)

- RCU transitions non-idle <-> idle:
  - kernel <-> user
  - kernel <-> idle
  - irq <-> user
  - irq <-> idle

**Each of these transitions cause rcu\_data::dynticks to be incremented.**

- RCU remains non-idle:
  - irq <-> nmi
  - nmi <-> irq
  - kernel <-> irq
  - Irq <-> kernel
- **Each of these transitions cause rcu\_data::dynticks to NOT be incremented.**

# nohz CPUs tick turn off : intro

Scheduler tick is turn off is attempted:

- From cpu idle loop: `tick_nohz_idle_stop_tick()`;
- At the end of IRQs:

```
irq_exit()  
    -> tick_nohz_irq_exit()  
        -> tick_nohz_full_update_tick()
```

`tick_nohz_full_update_tick()` checks `tick_dep_mask`. If any bits are set, tick is not turned off.

# nohz CPUs : reasons to keep tick ON

Introduced in commit

d027d45d8a17 ("nohz: New tick dependency mask")

- POSIX CPU timers
- Perf events
- Scheduler
  - (for nohz\_full, tick may be turned off if RQ has only 1 task)

# The problem: nohz\_full and RCU (TODO: Diagram)

Consider the scenario:

- A nohz\_full CPU with tick turned off is looping in kernel mode.
- The GP kthread on another CPU notices CPU is taking too long.
- Sets per-cpu hints: rcu\_urgent\_qs to true and sends resched\_cpu().
- IPI runs on nohz\_full CPU.
- It enters the scheduler, still nothing happens. Why??
- Grace period continues to stall... resulting in OOM.

Reason:

- Entering the scheduler only means the per-CPU quiescent state is updated
- Does not mean the global view of the CPU's quiescent state is updated...

# The problem: nohz\_full and RCU

What does updating global view of a CPU QS mean?

Either of the following:

- A. Quiescent state propagates up the RCU TREE to the root
- B. The `rcu_data :: dynticks` counter of the CPU is updated.
  - This causes the FQS loop to do A on dyntick transition of CPU.

Either of these operations will result in a globally updated view of CPU QS

Both of these are expensive and need to be done carefully.

# The problem: nohz\_full and RCU

The global view of a CPU QS is updated only from a few places..

- From the RCU softirq (which is raised from scheduler tick)
  - This does A (tick).
- Kernel/IRQ -> Idle / User transition
  - This does B (dyntick transition).
- cond\_resched() in PREEMPT=n kernels
  - This does B (dyntick transition).
- Other special cases where a GP must not be stalled and an RCU reader is not possible: (rcutorture, multi\_cpu\_stop() etc.)



# The problem: nohz\_full and RCU

Possible Solution:

Use option B (forced dynticks transition) in the IPI handler.

- This would be rather hackish. Currently we do this only in very special cases.
- This will only solve the issue for the current GP, future GPs may still suffer.

# The problem: nohz\_full and RCU

Final Solution: Turn back the tick on.

- Add a new “RCU” tick dependency mask.
- RCU FQS loop sets rcu\_urgent\_qs to true and does resched\_cpu() on the stalling CPU as explained.
- IPI handler arrives at CPU.
- On IPI entry, we set tick dependency mask.
- On IPI exit, the sched-tick code turns the tick back on.

This solution prevents the OOM splats.

<https://git.kernel.org/pub/scm/linux/kernel/git/paulmck/linux-rcu.git/commit/?h=dev&id=5b451e66c8f78f301c0c6eaacb0246eb13e09974>

<https://git.kernel.org/pub/scm/linux/kernel/git/paulmck/linux-rcu.git/commit/?h=dev&id=fd6d2b116411753fc929e6523d2ada227aa553>

# dyntick counter cleanup attempts

## Why?

In the previous solution, found a bug where dyntick\_nmi\_nesting was being set to (DYNTICK\_IRQ\_NONIDLE | 2)

However, the code was checking (dyntick\_nmi\_nesting == 2)

## Link to buggy commit:

<https://git.kernel.org/pub/scm/linux/kernel/git/paulmck/linux-rcu.git/commit/?h=dev&id=13c4b07593977d9288e5d0c21c89d9ba27e2ea1f>

**However, Paul wants the existing mechanism redesigned** rather than cleaned up since it's too complex: <http://lore.kernel.org/r/20190828202330.GS26530@linux.ibm.com>

# Thank you ...

- For questions, please email the list: [rcu@vger.kernel.org](mailto:rcu@vger.kernel.org)
- Follow us on Twitter:
  - paulmckrcu@ joel\_linux@ boqun\_feng@ srostedt@