

Workqueue insights

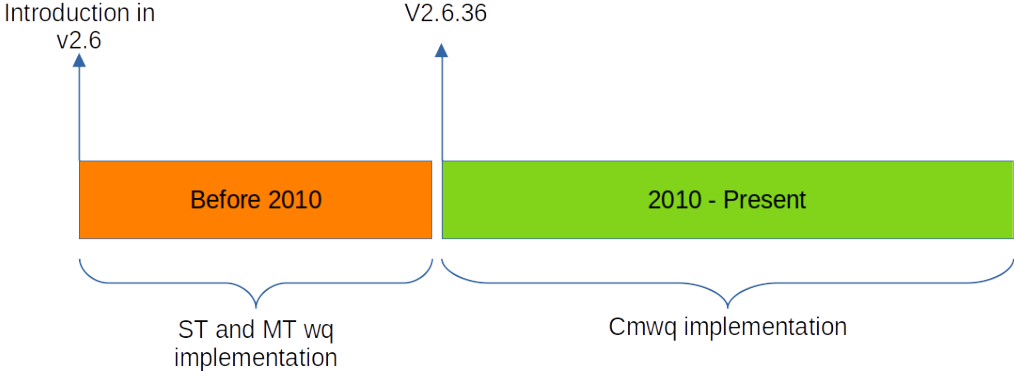
Prathu Baronia & Neeraj Upadhyay

29-07-2023

Before we begin

- ▶ wq: workqueue
- ▶ cmwq: concurrency managed workqueue
- ▶ Kernel version considered: v6.5-rc3

History



WQ API

Used for asynchronous execution of functions in process context

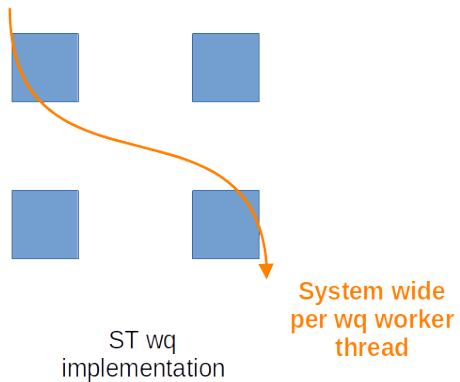
API

- `alloc_workqueue()`
- `struct workqueue_struct`
- `queue_work|on()`
- `schedule_work|on()`

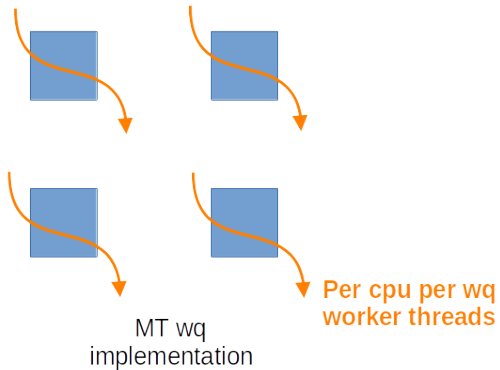
Remember

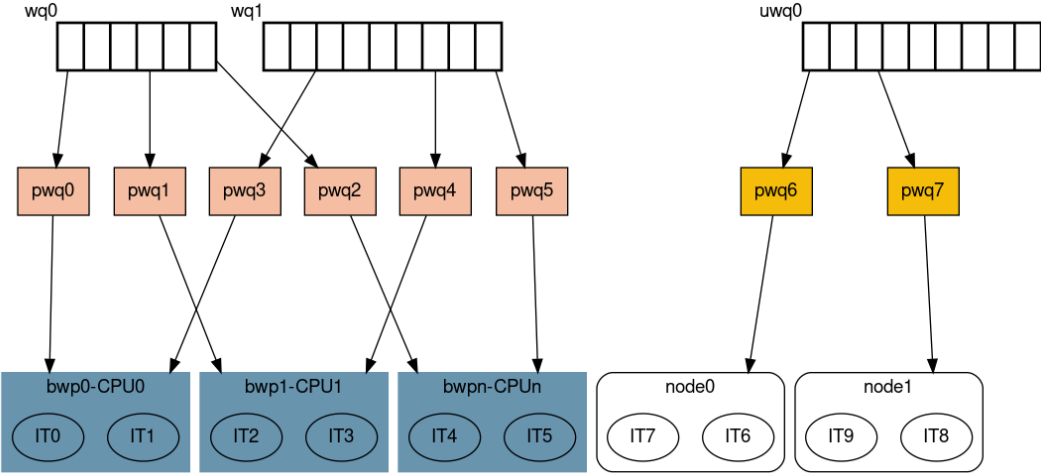
- Execution thread can sleep
- No userspace memory access
- Self queuing functions are allowed

Legacy implementation: ST



Legacy implementation: MT





Workqueue software layers

Workqueue

Bounded norm/high prio

Unbounded norm/high prio

Freezable

Ordered

Power Efficient

Mem Reclaim

Attributes:

- max_active
- Unbound - cpu_mask, no_numa, nice

Pool Workqueue

Bounded

CPU0 CPU1 ... CPU_n

Per CPU

Unbounded

Default PWQ

NUMA Node PWQ

Node 0	PWQ _x
Node 1	PWQ _y

- Flush management
- In flight works
- Active/Inactive works

Worker Pool

Bounded norm/high prio

Unbounded norm/high prio default

NUMA node unbound Pools

Custom Unbound pools - cpumask, nice

Worklist

Idle List

Busy Hash

Workers

Bounded

kworker/0:0

kworker/1:0H

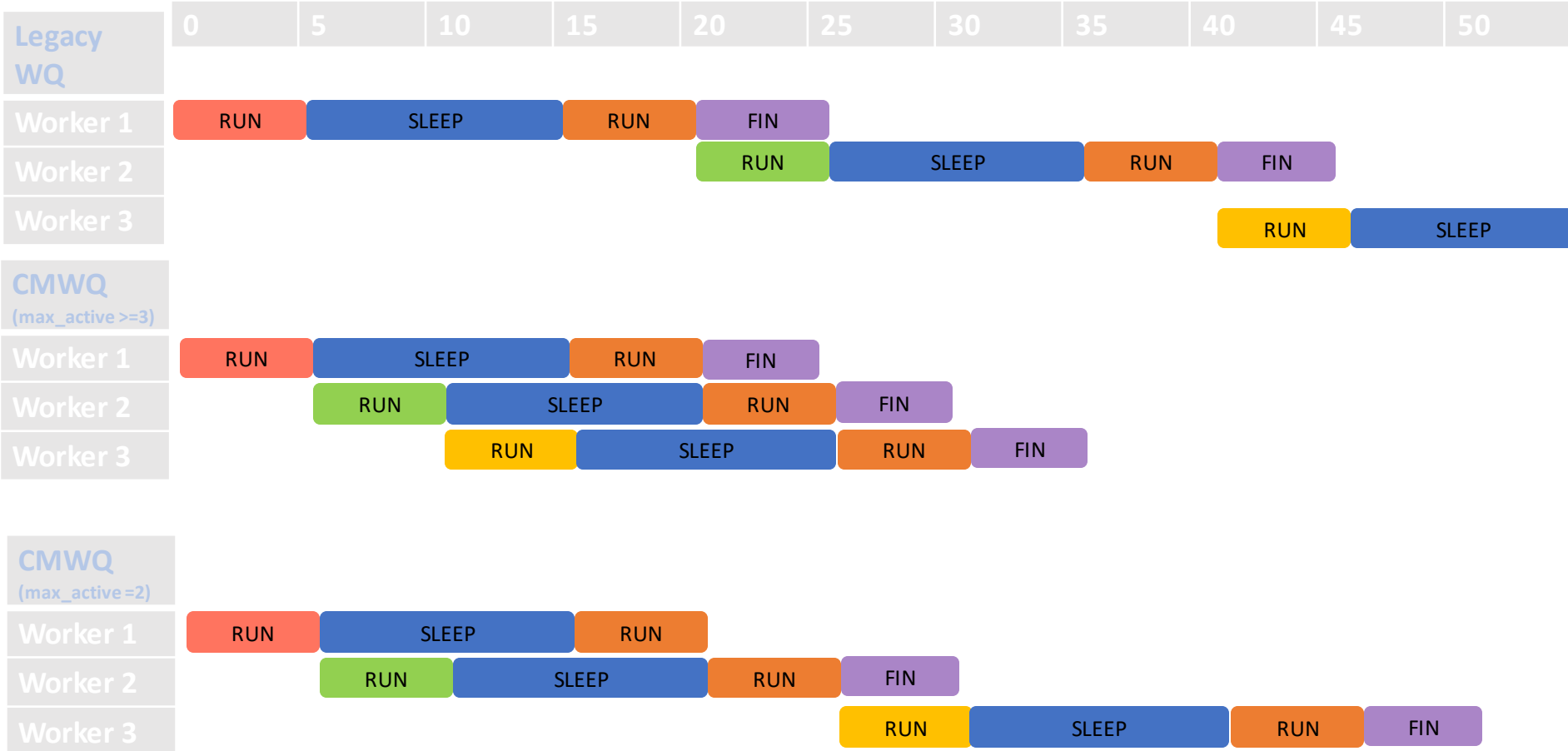
Per CPU/Concurrency Managed, unless executing CPU intensive work

Unbounded

kworker/u3:4

No Concurrency Management, unbounded/NUMA node CPU affinity.

WQ EXECUTION SCENARIOS



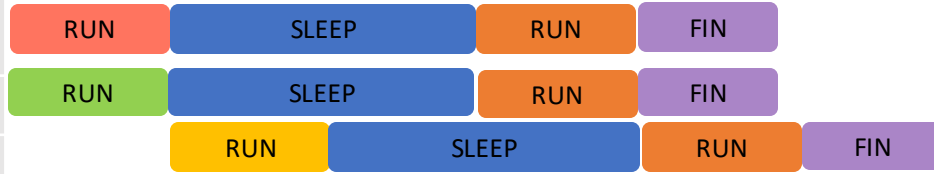


CMWQ
(max_active >=3)

Worker 1
(CPU Intensive)

Worker 2

Worker 3



CMWQ
(max_active >=3)

Worker 1

Worker 2
(CPU Intensive)

Worker 3

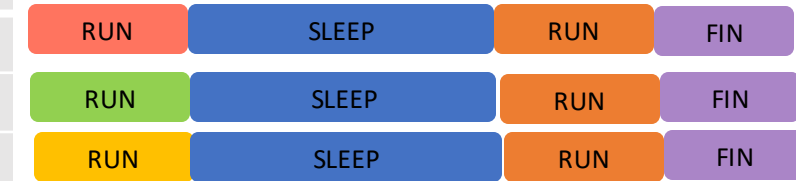


CMWQ
(Unbounded)

Worker 1

Worker 2

Worker 3



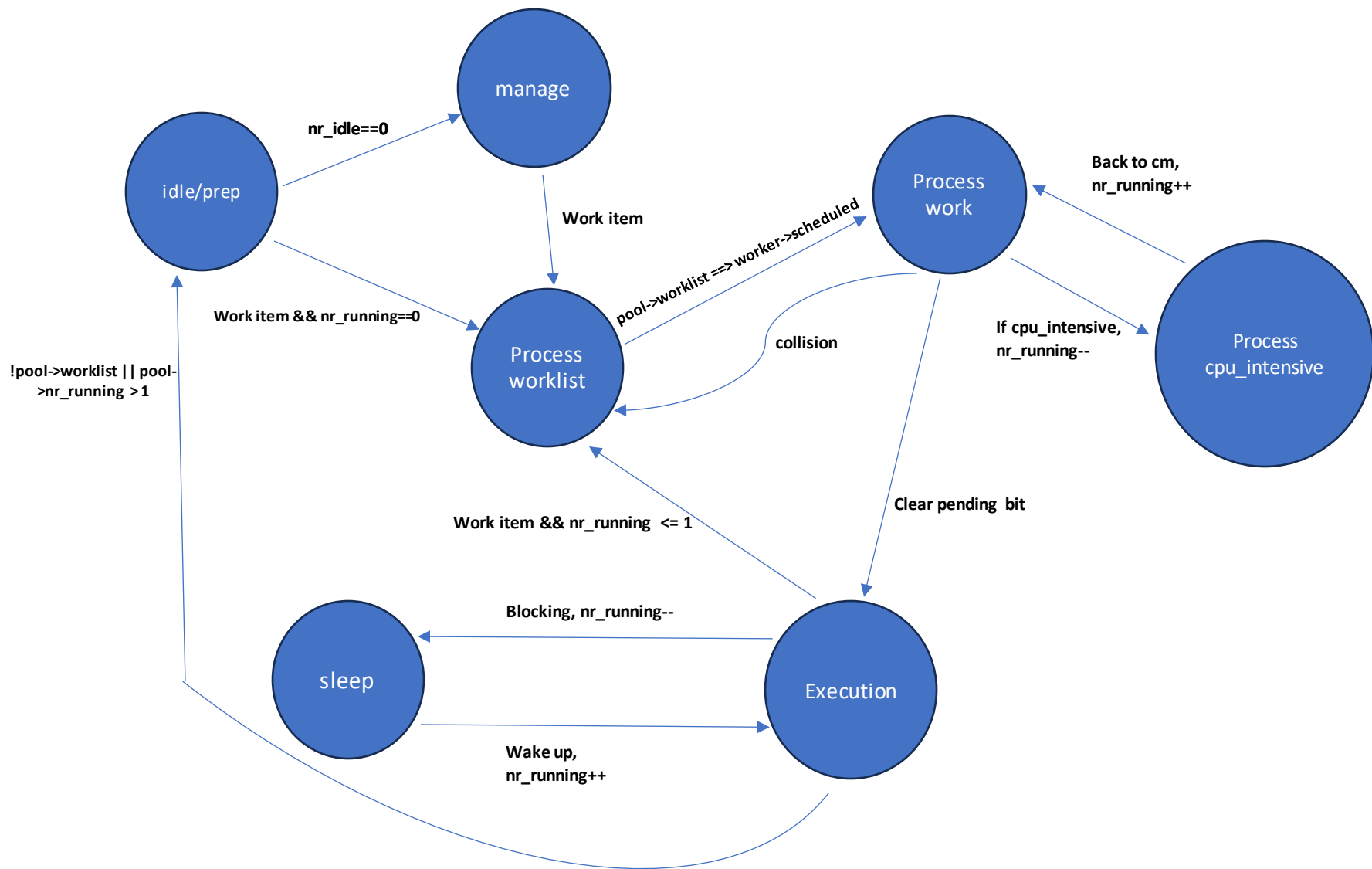
WQ EXECUTION SCENARIOS
CONTD.



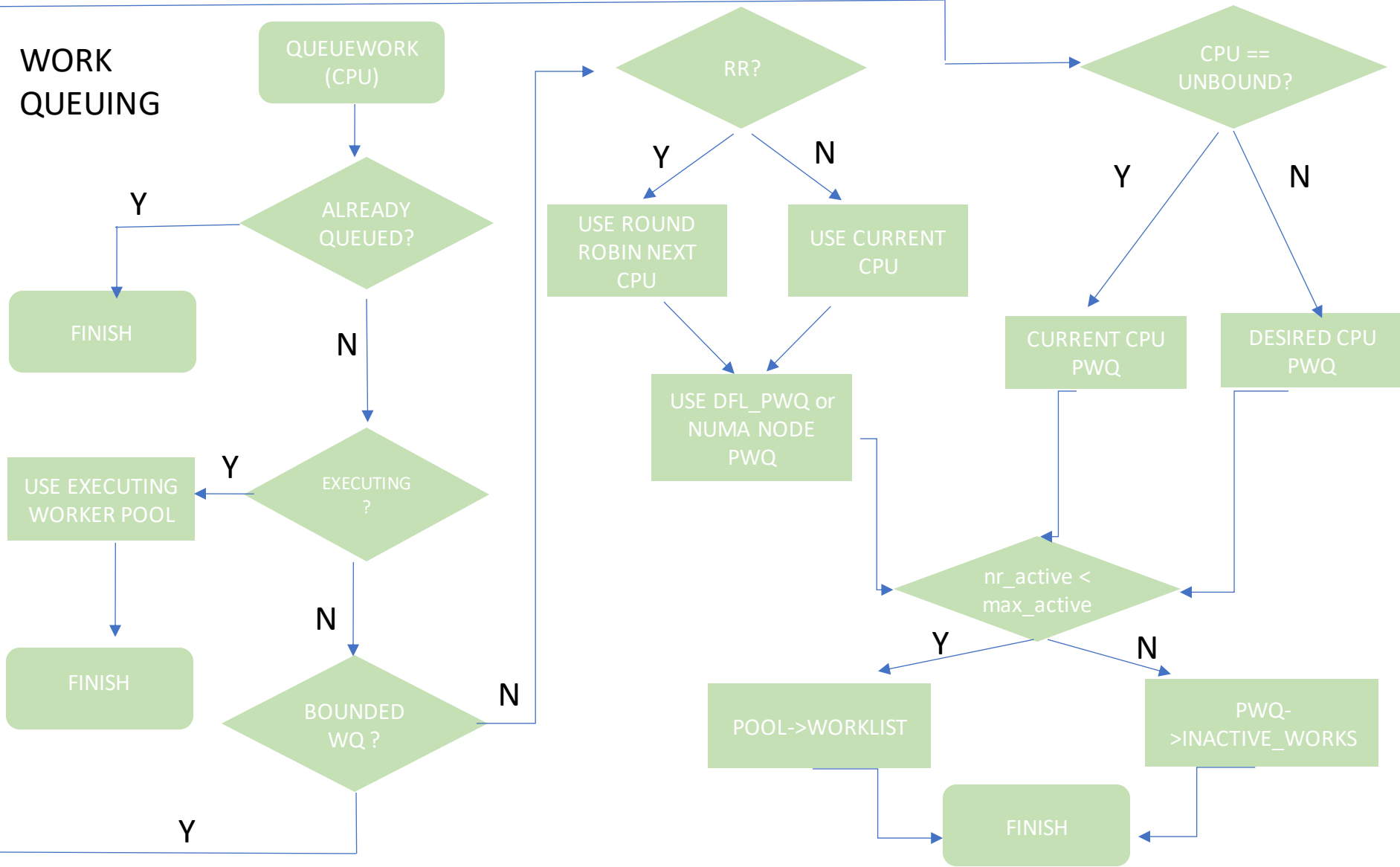
CMWQ (Ordered)
Worker 1
Worker 2
Worker 3



WQ EXECUTION SCENARIOS
CONTD.



WORK QUEUING



System wide wqs

- ▶ `system_wq`
- ▶ `system_highpri_wq`
- ▶ `system_long_wq`
- ▶ `system_unbound_wq`
- ▶ `system_freezable_wq`
- ▶ `system_power_efficient_wq`
- ▶ `system_freezable_power_efficient_wq`

FAQs

- ▶ Which system wq to use in my code?
- ▶ When to use unbounded wqs?
- ▶ How to decide `max_active`?

Usecases of system wide wqs

- ▶ `system_wq`:
 - ▶ `kernel/smp.c`: `smp_call_on_cpu()`
`queue_work_on(cpu, system_wq, &sscs.work);`
- ▶ `system_highpri_wq`:
 - ▶ `drivers/gpu/drm/radeon/radeon_display.c`
`radeon_x->flip_queue = alloc_workqueue("radeon-x", WQ_HIGHPRI, 0);`
- ▶ `system_long_wq`:
 - ▶ `drivers/base/core.c`: `devlink_dev_release()`
`queue_work(system_long_wq, &link->rm_work);`

Usecases of system wide wqs cont.

- ▶ `system_unbound_wq`:
 - ▶ `drivers/base/dd.c`: `driver_deferred_probe_trigger()`
`queue_work(system_unbound_wq, &deferred_probe_work);`
- ▶ `system_freezable_wq`:
 - ▶ `drivers/virtio/virtio_balloon.c`: `update_balloon_size_func()`
`queue_work(system_freezable_wq, work);`
- ▶ `system_power_efficient_wq`:
 - ▶ `sound/soc/codecs/rt5645.c`: `rt5645_irq()` irq handler
`queue_delayed_work(system_power_efficient_wq,`
`&rt5645->jack_detect_work, msecs_to_jiffies(250));`