

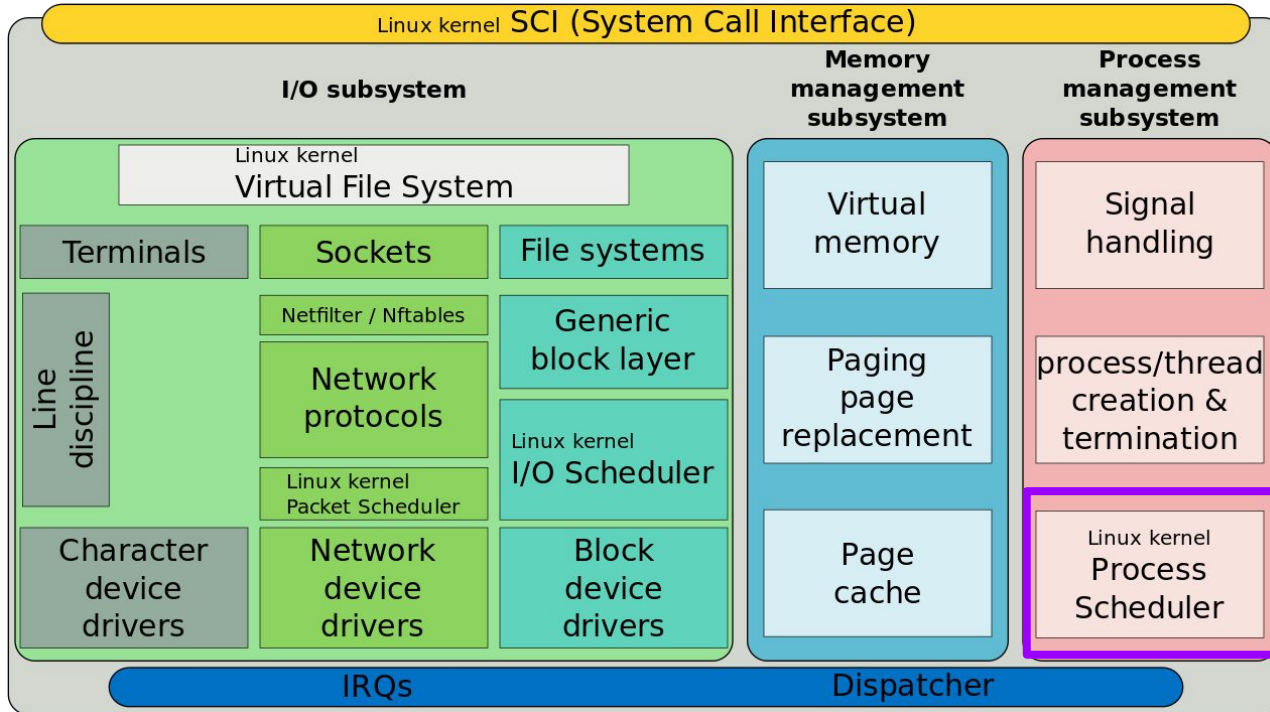
ELEC 424/553

Mobile & Embedded Systems

Lecture 6 - How Does Linux Schedule Processes? (2)

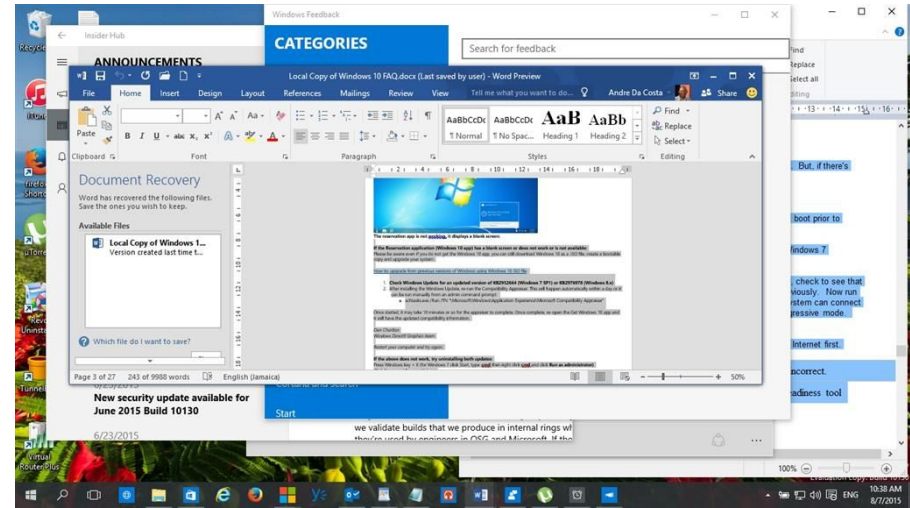
The saga continues...

Where Are We In The Linux Kernel?



Scheduling Processes

- Goal: Run a bunch of applications
- **Multitasking** makes it look like OS is simultaneously executing more than 1 application/process
 - Capable of alternating between processes



Andre Da Costa, "How to: manage running programs and virtual desktops using Task View in Windows 10". URL: <https://answers.microsoft.com/en-us/insider/forum/all/how-to-manage-running-programs-and-virtual/17d068b7-5e4a-4351-a019-afa528a81538>

Two Modern Scheduling Concerns

1. Priority
2. Timeslice

Priority & Nature of Scheduler Classes

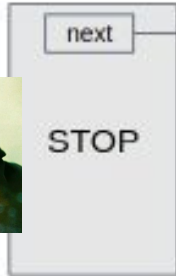
Priority & Nature of Scheduler Classes

Highest Priority  Lowest Priority

Priority & Nature of Scheduler Classes

Literally halts
everything
Unstoppable

Process
migration



Highest Priority



Lowest Priority

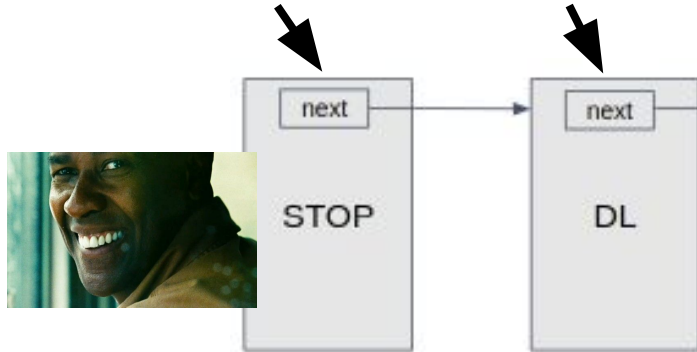
Priority & Nature of Scheduler Classes

Literally halts
everything
Unstoppable

Process
migration

Hard deadline
processes
soonest first

Video



Highest Priority

Lowest Priority

Priority & Nature of Scheduler Classes

Literally halts
everything
Unstoppable

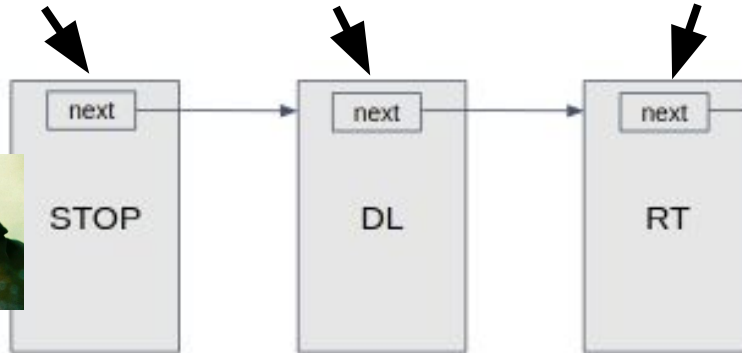
Process
migration

Hard deadline
processes
soonest first

Video

Realtime
Latency-important
processes

Interrupt request
(IRQ)

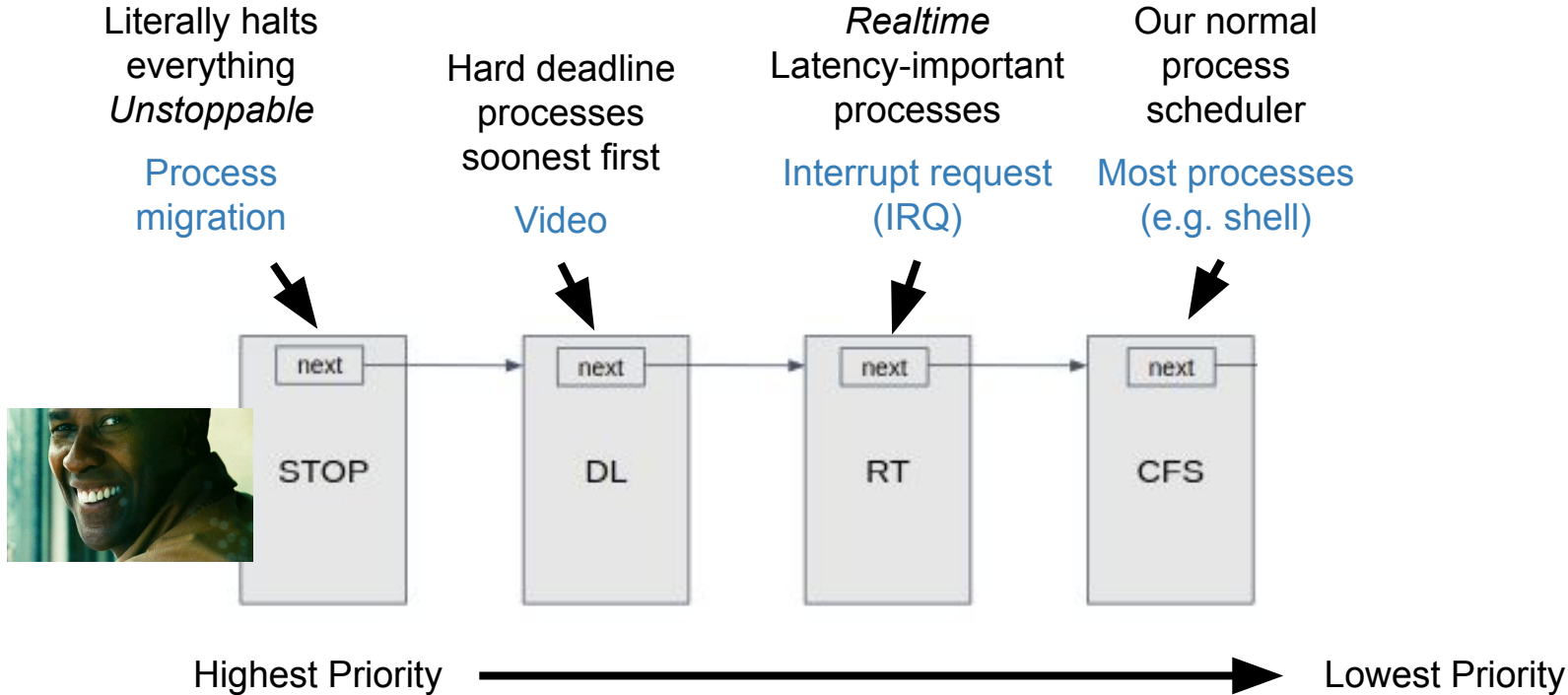


Highest Priority

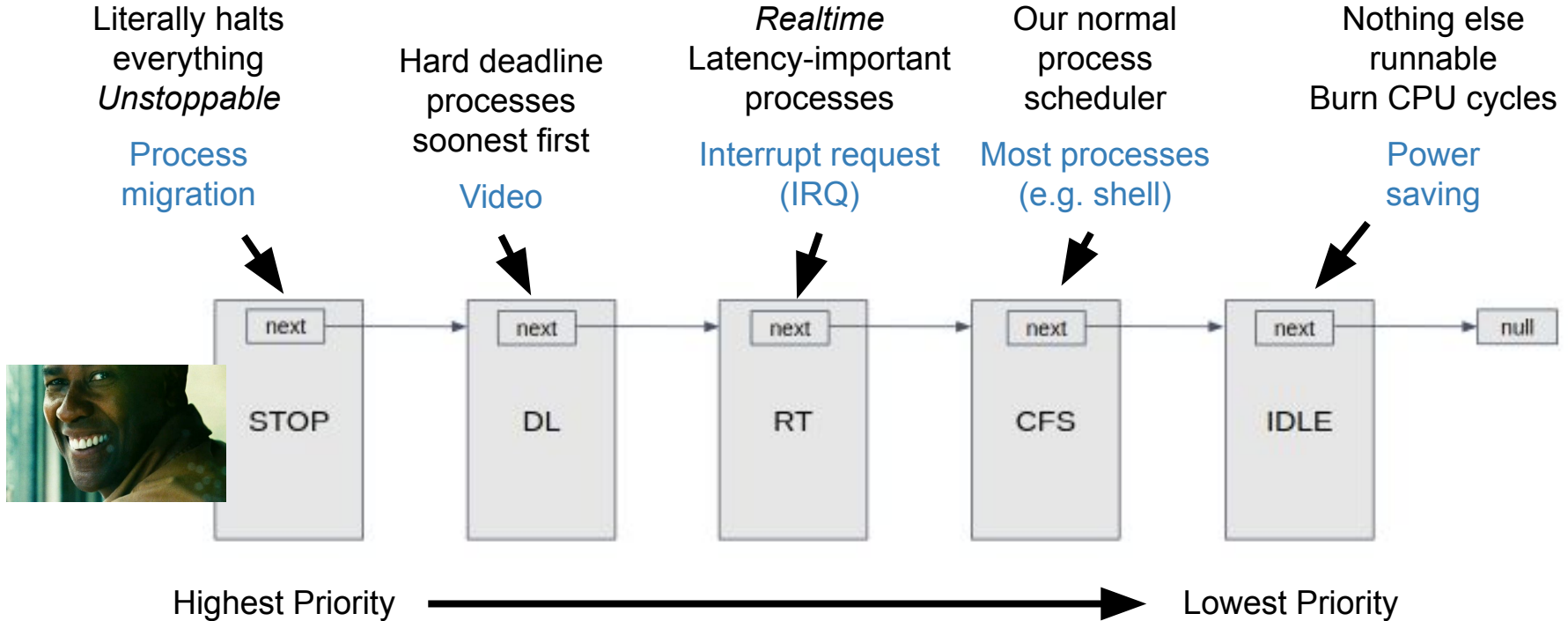
Lowest Priority



Priority & Nature of Scheduler Classes



Priority & Nature of Scheduler Classes



Scheduler Classes

```
/ kernel / sched / core.c All
5902  * WARNING: must be called with preemption disabled!
5903  */
5904  static void __sched notrace __schedule(bool preempt)
5905  {
5906      struct task_struct *prev, *next;
5907      unsigned long *switch_count;
5908      unsigned long prev_state;
5909      struct rq_flags rf;
5910      struct rq *rq;
5911      int cpu;
5912
5913      cpu = smp_processor_id();
5914      rq = cpu_rq(cpu);
5915      prev = rq->curr;
5916
```

```
/ kernel / sched / sched.h All sym Search
904
905  /*
906   * This is the main, per-CPU runqueue data structure.
907   *
908   * Locking rule: those places that want to lock multiple runqueues
909   * (such as the load balancing or the thread migration code), lock
910   * acquire operations must be ordered by ascending &runqueue.
911   */
912  struct rq {
913      /* runqueue lock: */
914      raw_spinlock_t _lock;
915
916      /*
917       * nr_running and cpu_load should be in the same cacheline because
918       * remote CPUs use both these fields when doing load calculation.
919       */
920      unsigned int uclamp_flags;
921
922      #define UCLAMP_FLAG_IDLE 0x01
923      #endif
924
925      struct cfs_rq cfs;
926      struct rt_rq rt;
927      struct dl_rq dl;
928
```

Linux 5.14.2 Source via Bootlin Elixir Cross Referencer

Left: <https://elixir.bootlin.com/linux/v5.14.2/source/kernel/sched/core.c#L5904>

Right: <https://elixir.bootlin.com/linux/v5.14.2/source/kernel/sched/sched.h#L912>

Kernel Calls Scheduler via `schedule()`

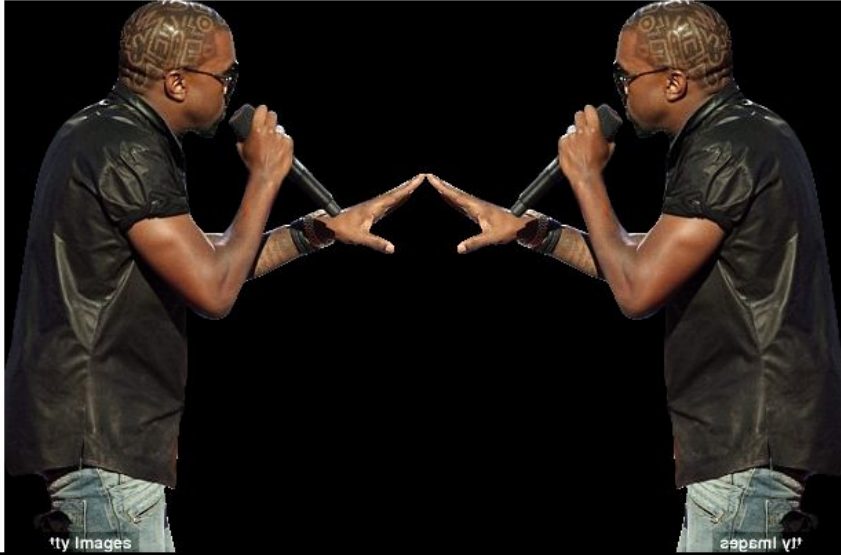
- Independent of scheduler classes
- Asks scheduler of class with greatest priority for process that should be executed
 - (if class has process ready to execute)

```
/ kernel / sched / core.c All sy▼  
6098  asmlinkage __visible void __sched schedule(void)  
6099  {  
6100      struct task_struct *tsk = current;  
6101  
6102      sched_submit_work(tsk);  
6103      do {  
6104          preempt_disable();  
6105          __schedule(false);  
6106          sched_preempt_enable_no_resched();  
6107      } while (need_resched());  
6108      sched_update_worker(tsk);  
6109  }  
6110  EXPORT_SYMBOL(schedule);
```

Linux Source via Bootlin Elixir Cross Referencer

<https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L6098>

Kernel Calls Scheduler via `schedule()`



We want to prevent this from happening

<https://knowyourmeme.com/photos/292972-kanye-interrupts-imma-let-you-finish>

```
sched / core.c All sy▼  
  
asmlinkage __visible void __sched schedule(void)  
{  
    struct task_struct *tsk = current;  
  
    sched_submit_work(tsk);  
    do {  
        preempt_disable();  
        __schedule(false);  
        sched_preempt_enable_no_resched();  
    } while (need_resched());  
    sched_update_worker(tsk);  
  
EXPORT_SYMBOL(schedule);
```

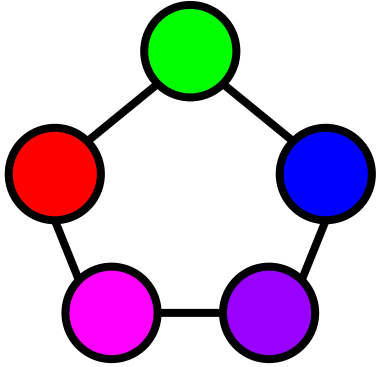
Linux Source via Bootlin Elixir Cross Referencer

<https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L6098>

A Brief History of Linux Schedulers

A Brief History of Linux Schedulers

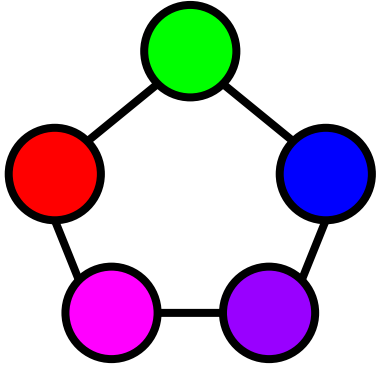
1.2: Circular Queue



Round-robin
Simple, fast

A Brief History of Linux Schedulers

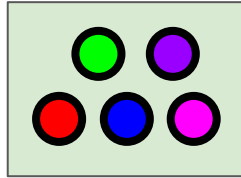
1.2: Circular Queue



Round-robin
Simple, fast

2.4: $O(n)$

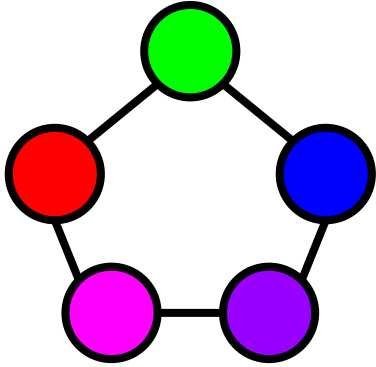
Which next?



Check all processes
Scalability issues

A Brief History of Linux Schedulers

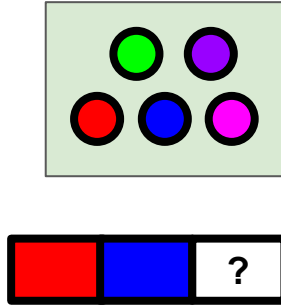
1.2: Circular Queue



Round-robin
Simple, fast

2.4: $O(n)$

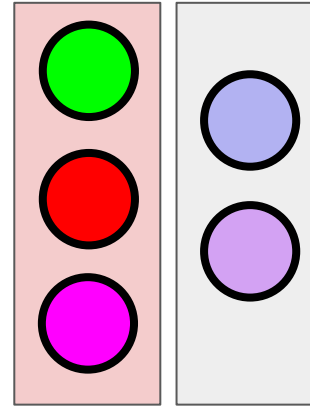
Which next?



Check all processes
Scalability issues

2.6: $O(1)$

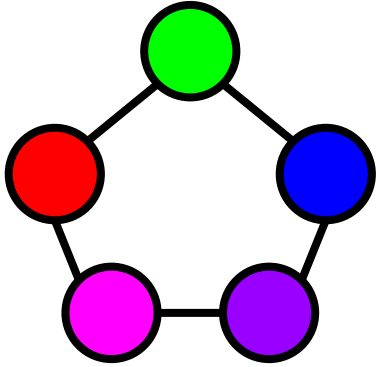
Active Expired



Constant time
Arrays

A Brief History of Linux Schedulers

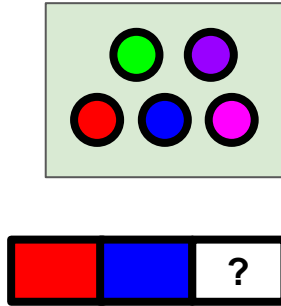
1.2: Circular Queue



Round-robin
Simple, fast

2.4: $O(n)$

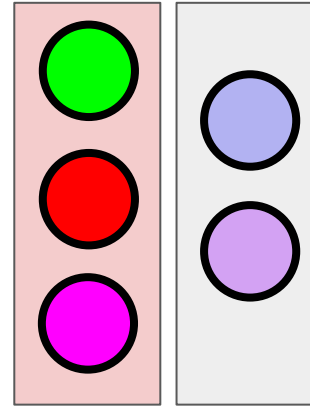
Which next?



Check all processes
Scalability issues

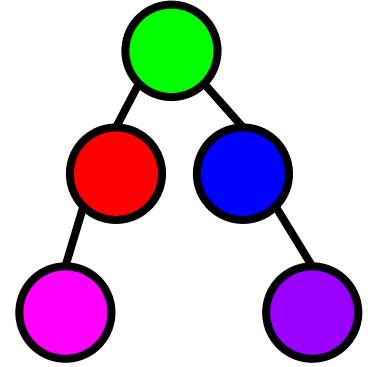
2.6: $O(1)$

Active Expired



Constant time
Arrays

2.6.23: CFS



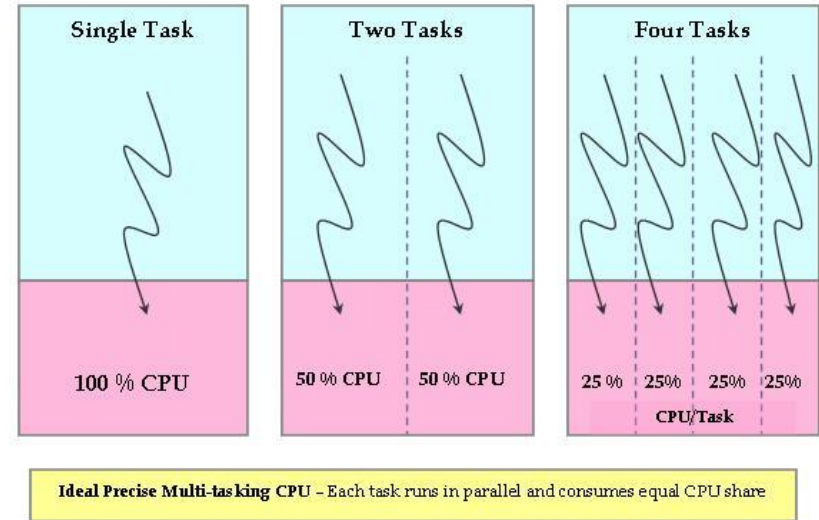
Red-black tree
Virtual runtime

CFS Has Arrived



The Idea Behind Linux's Completely Fair Scheduler (CFS)

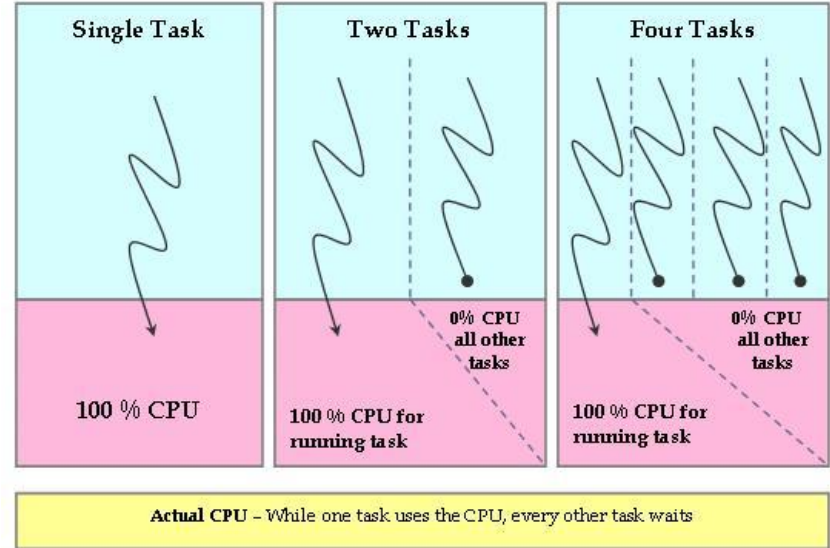
- Just imagine a perfect processor!
- n processes that can execute
- Split processor's schedule into $1/n$ length blocks
- Drive schedule down to infinitesimal length of time
- Could say that processes were executed using $(1/n) \times 100\%$ of the processor's power
- “Perfect multitasking”



Chandandeep Singh Pabla, “Completely Fair Scheduler”. URL: <https://www.linuxjournal.com/node/10267>

Reality Sinks In

- Processor can only run one process at a time!
- Rapidly alternating processes (preempting) would be terrible
- Instead, set a schedule period and give equal “virtual time” to processes
- By *simply* redefining time, we become *fair*



Chandandeep Singh Pabla, "Completely Fair Scheduler". URL:
<https://www.linuxjournal.com/node/10267>

Math Behind CFS: Target Latency & Timeslice

- **Target latency:** Total amount of time established by CFS to allow each task to get a turn

- `sysctl_sched_latency` (ms)

- Processes given weights based on `nice` (see table)

nice	-20	-19	-18	-17	-16	-15	-14	-13	-12	-11
weight	88761	71755	56483	46273	36291	29154	23254	18705	14949	11916
nice	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1
weight	9548	7620	6100	4904	3906	3121	2501	1991	1586	1277
nice	0	1	2	3	4	5	6	7	8	9
weight	1024	820	655	526	423	335	272	215	172	137
nice	10	11	12	13	14	15	16	17	18	19
weight	110	87	70	56	45	36	29	23	18	15

Nice-to-weight conversion

Jinkyu Koo, "Linux kernel scheduler". URL: <https://helix979.github.io/jkoo/post/os-scheduler/>

- $$timeslice = (target\ latency) \times \frac{(process\ weight)}{(sum\ of\ all\ process\ weights)}$$

Math Behind CFS (2): Virtual Runtime

- **Virtual runtime**, stored as `vruntime`, is warped version of physical runtime (wall-time)
- Warping based on weight (priority)
- Higher priority (`nice < 0`):
 - $vruntime < \text{physical runtime}$
- Lower priority (`nice > 0`):
 - $vruntime > \text{physical runtime}$

nice	-20	-19	-18	-17	-16	-15	-14	-13	-12	-11
weight	88761	71755	56483	46273	36291	29154	23254	18705	14949	11916
nice	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1
weight	9548	7620	6100	4904	3906	3121	2501	1991	1586	1277
nice	0	1	2	3	4	5	6	7	8	9
weight	1024	820	655	526	423	335	272	215	172	137
nice	10	11	12	13	14	15	16	17	18	19
weight	110	87	70	56	45	36	29	23	18	15

Nice-to-weight conversion

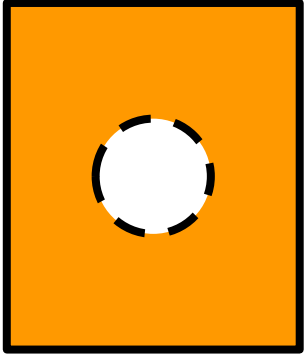
Jinkyu Koo, "Linux kernel scheduler". URL: <https://helix979.github.io/jkoo/post/os-scheduler/>

$$vruntime = (\text{physical runtime}) \cdot \frac{1024}{weight}$$

CFS Example

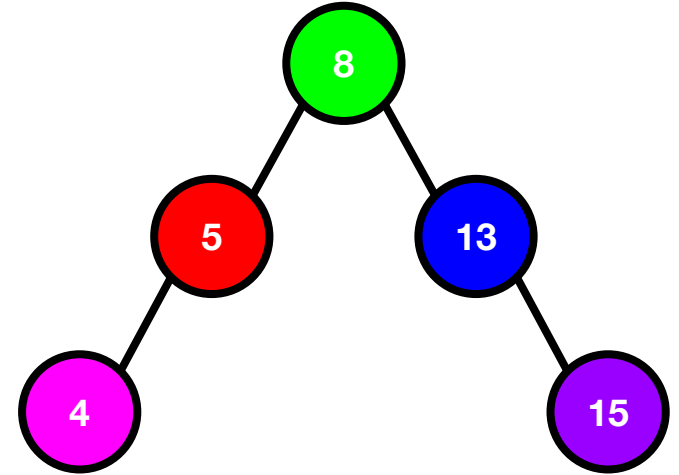
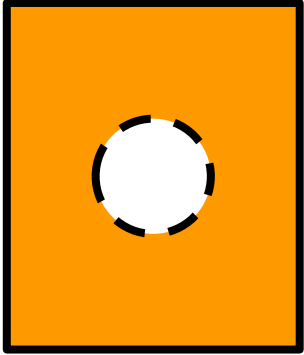
CFS Example

Processor



CFS Example

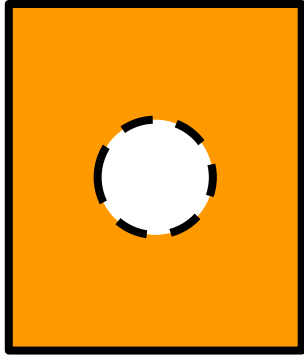
Processor



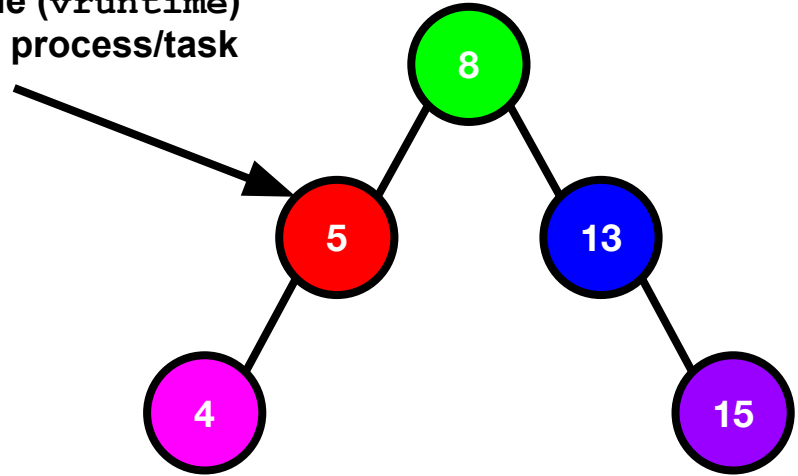
CFS Run Queue
Red-Black Tree

CFS Example

Processor



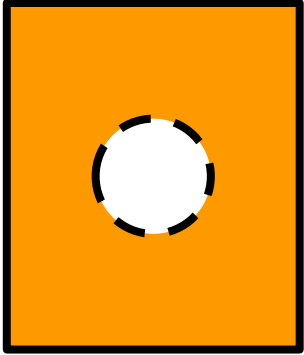
Virtual Runtime (vruntime)
for a runnable process/task



CFS Run Queue
Red-Black Tree

CFS Example

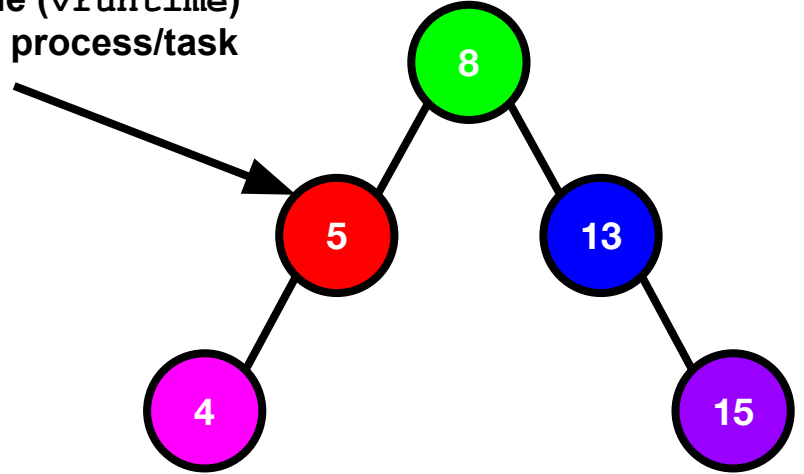
Processor



Processor Virtual Runtime Timeline



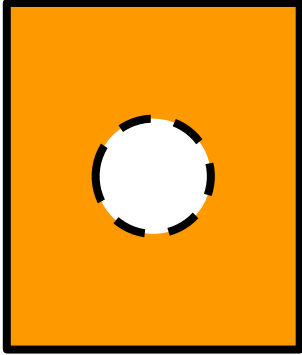
Virtual Runtime (vruntime)
for a runnable process/task



CFS Run Queue
Red-Black Tree

CFS Example

Processor



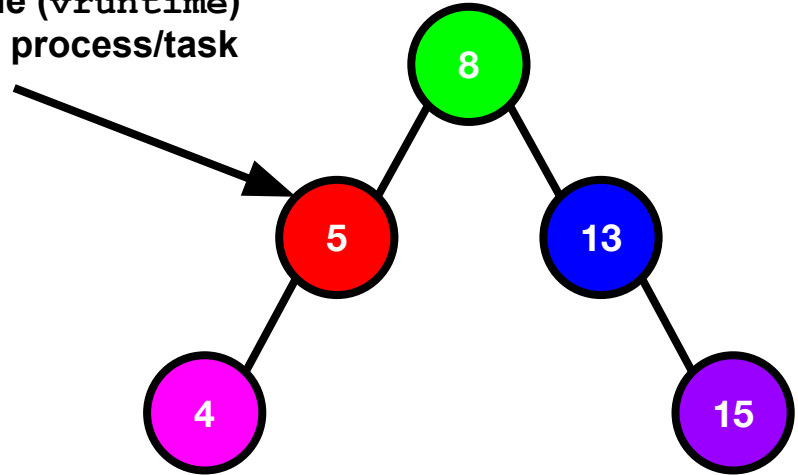
Processor Virtual Runtime Timeline



Processor Wall-Time Timeline

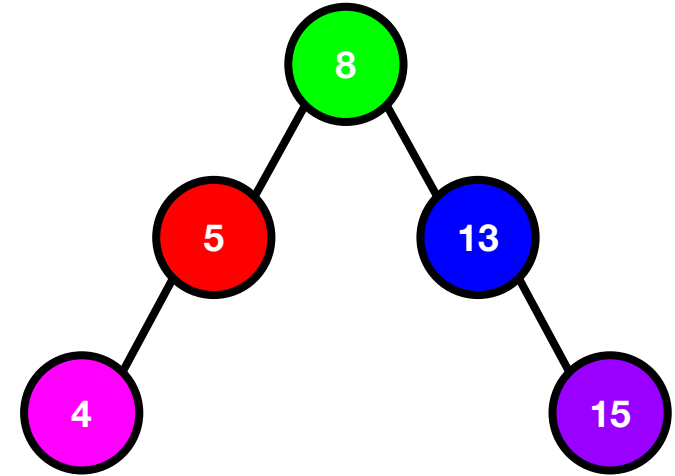
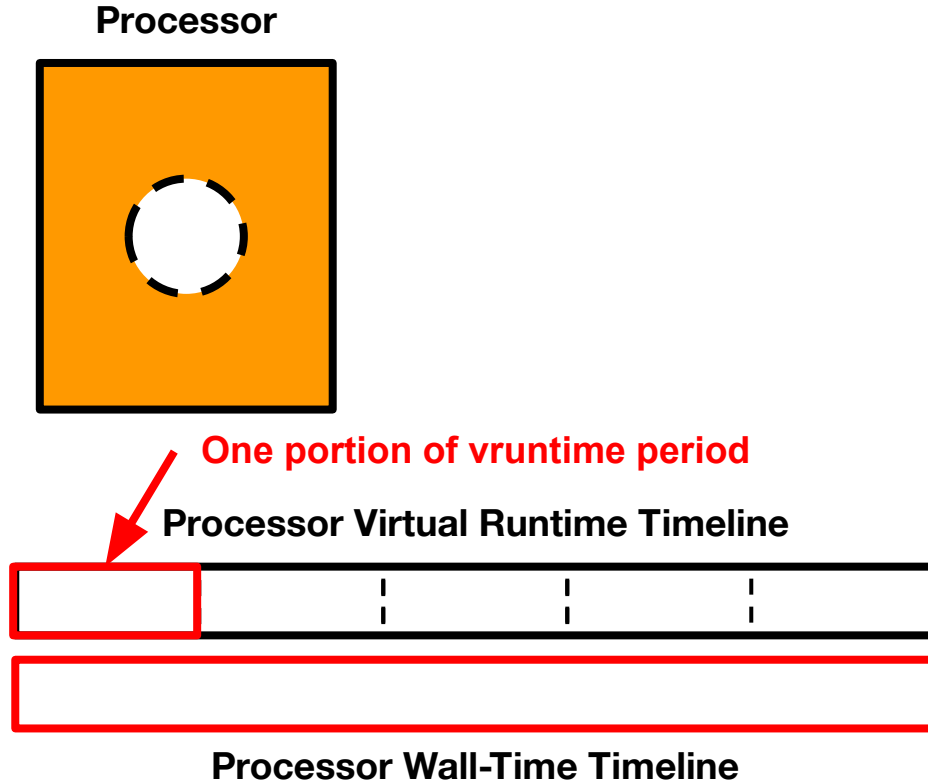


Virtual Runtime (vruntime)
for a runnable process/task



CFS Run Queue
Red-Black Tree

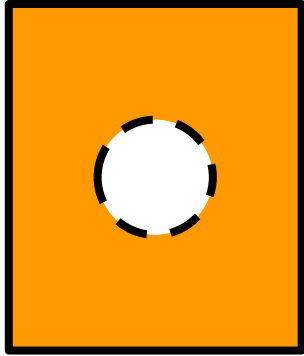
CFS Example



CFS Run Queue
Red-Black Tree

CFS Example

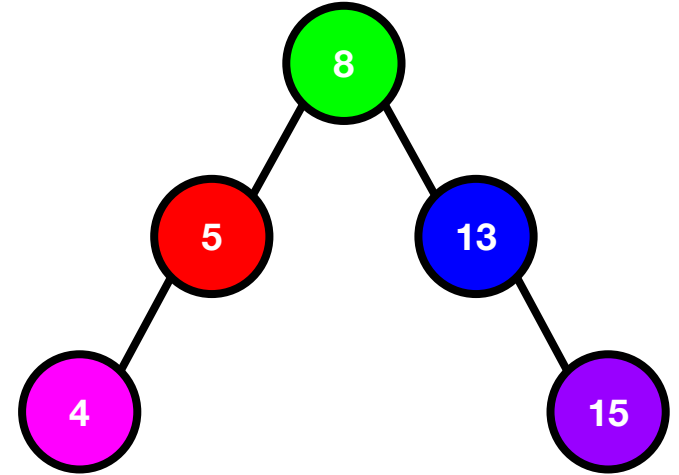
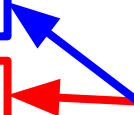
Processor



Processor Virtual Runtime Timeline



Processor Wall-Time Timeline

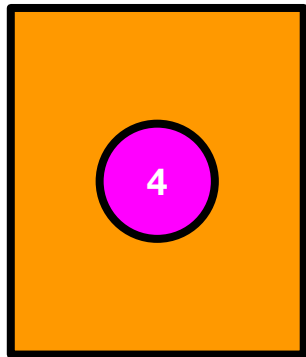


CFS Run Queue
Red-Black Tree

Can be different lengths of time

CFS Example

Processor



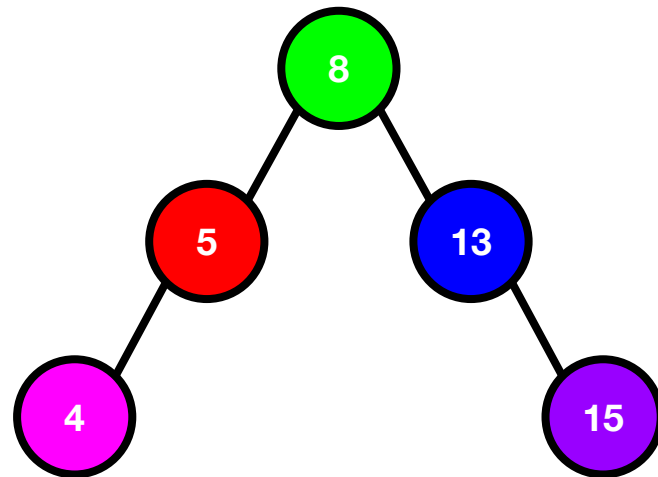
We chose this process because it has the lowest cumulative
vruntime

An advantage of a RB tree is that the process with the lowest vruntime will always be to the bottom left

Processor Virtual Runtime Timeline



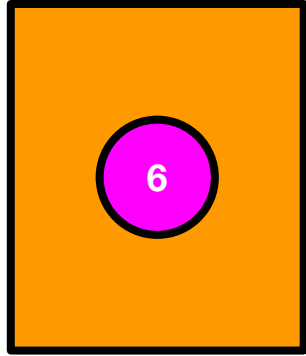
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Processor

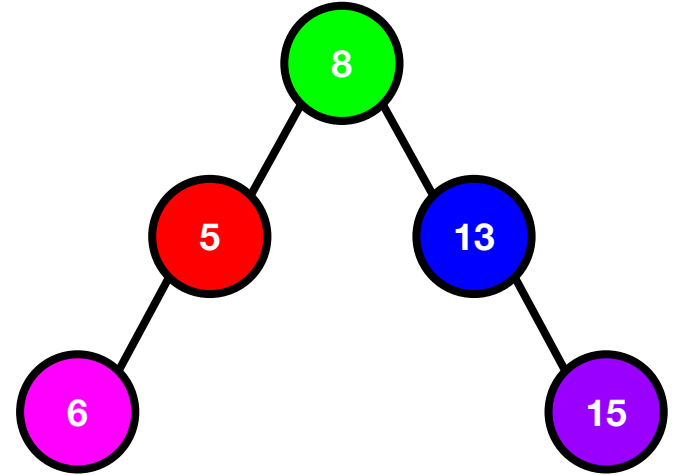


Less nice (nice<0)
Higher priority

Processor Virtual Runtime Timeline



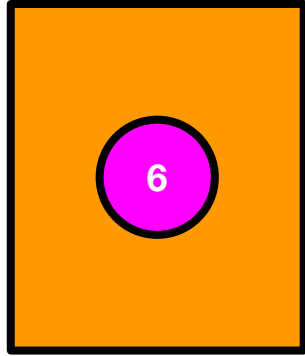
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Processor



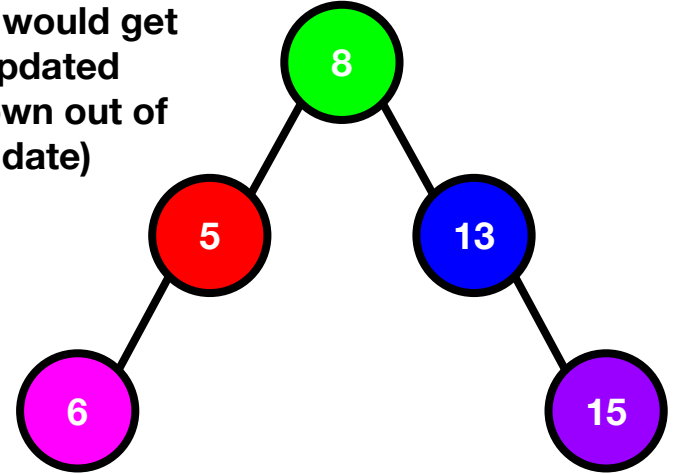
Less nice (nice<0)
Higher priority

Processor Virtual Runtime Timeline



Processor Wall-Time Timeline

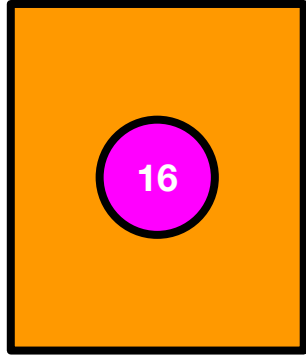
Tree would get
updated
(shown out of
date)



CFS Run Queue
Red-Black Tree

CFS Example

Processor

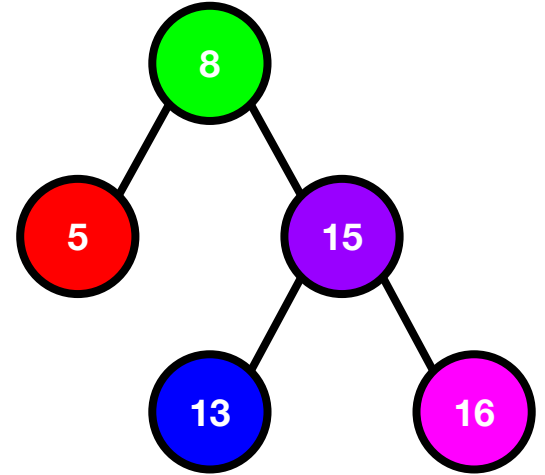


Less nice (nice<0)
Higher priority

Processor Virtual Runtime Timeline



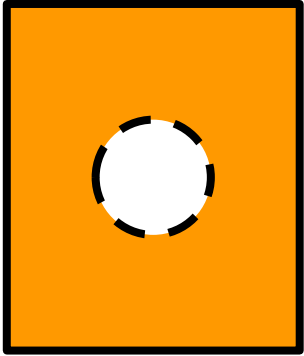
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

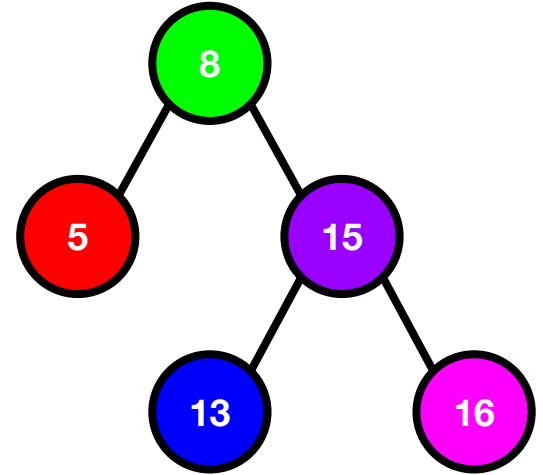
Processor



Processor Virtual Runtime Timeline

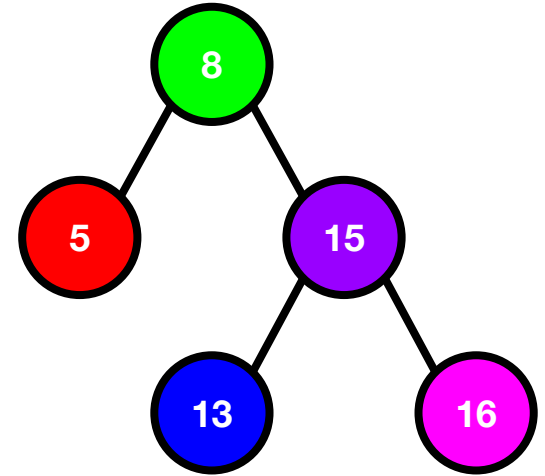
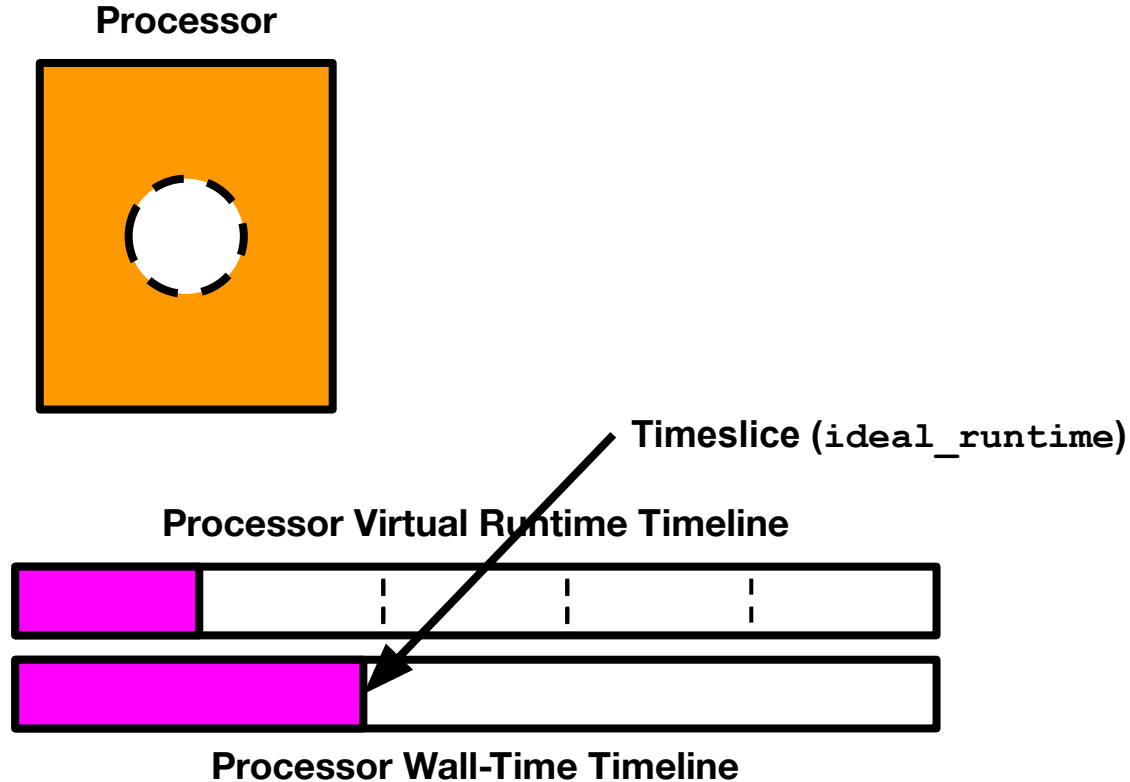


Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

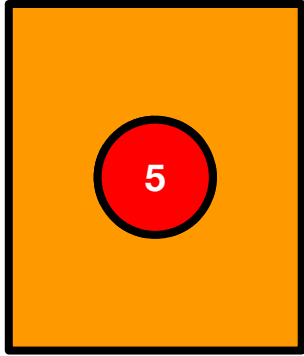
CFS Example



CFS Run Queue
Red-Black Tree

CFS Example

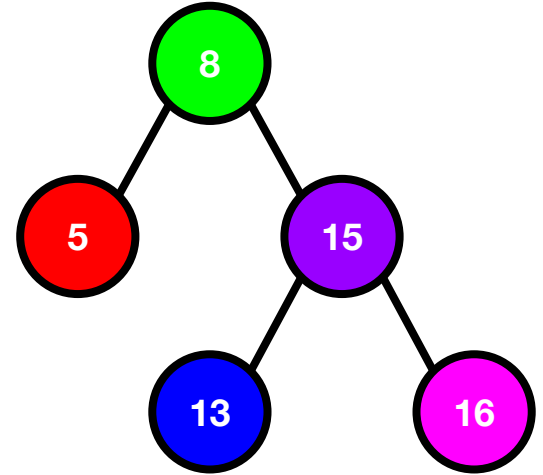
Processor



Processor Virtual Runtime Timeline



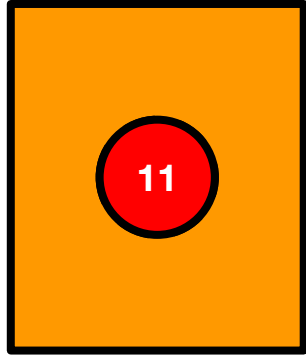
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Processor

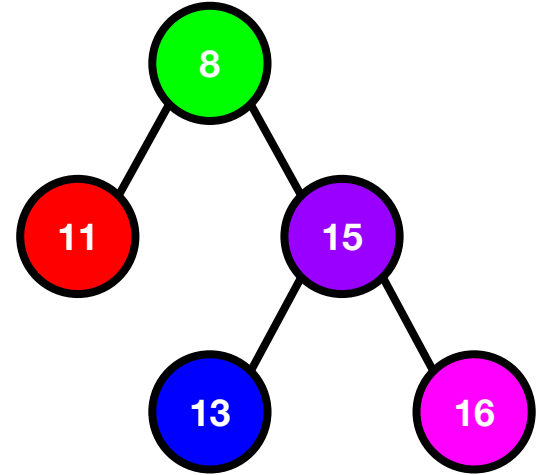


Default nice (nice=0)
Normal priority

Processor Virtual Runtime Timeline

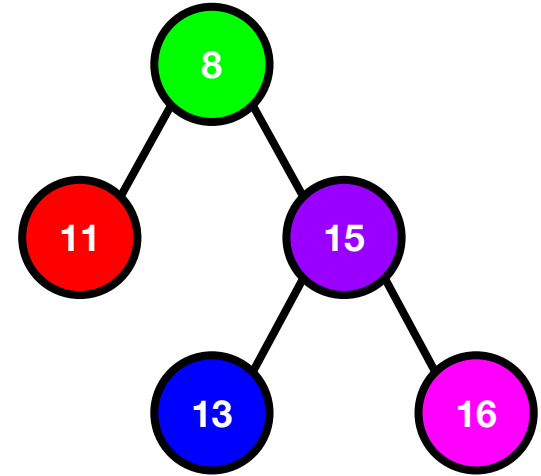
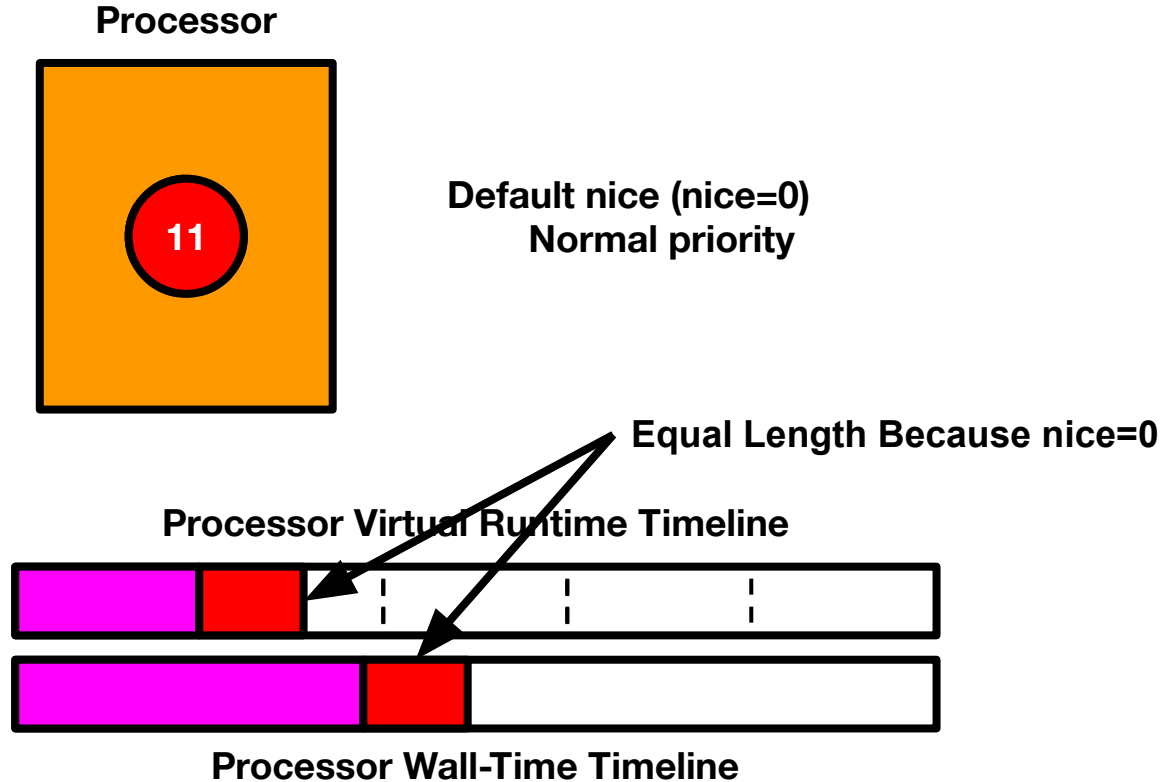


Processor Wall-Time Timeline



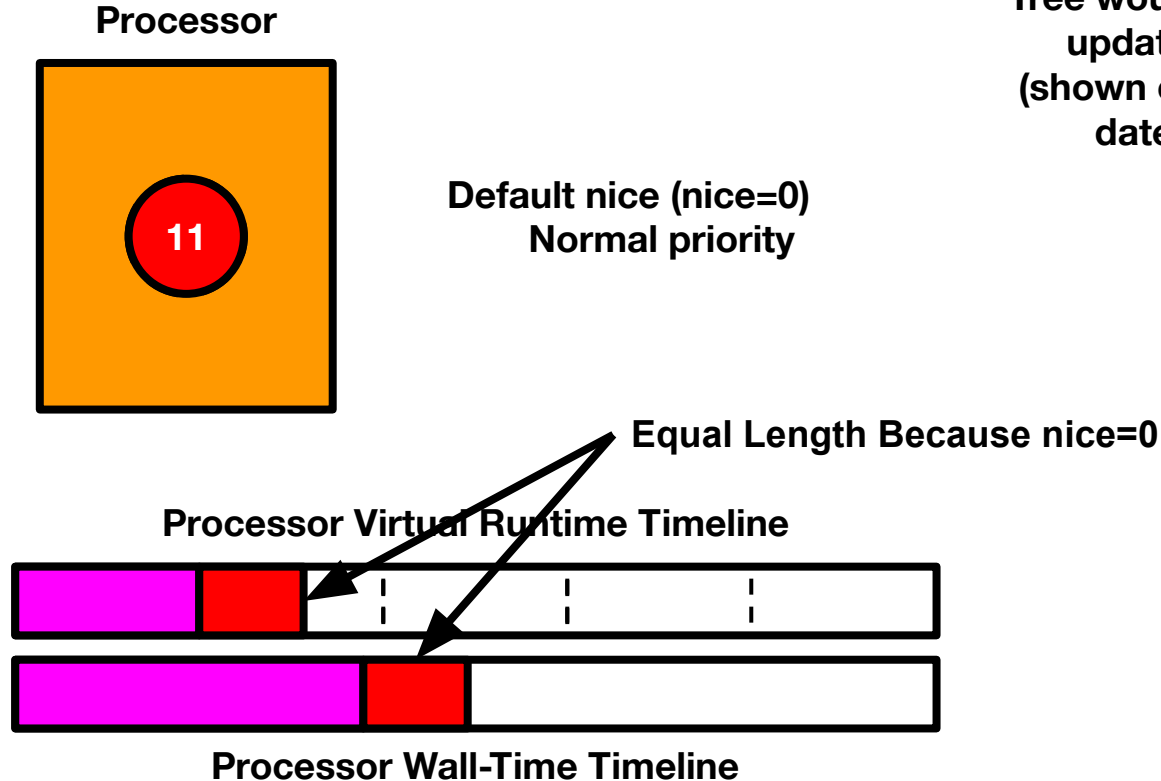
CFS Run Queue
Red-Black Tree

CFS Example

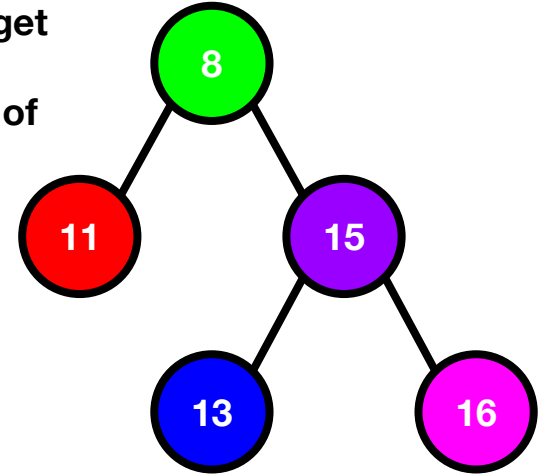


CFS Run Queue
Red-Black Tree

CFS Example



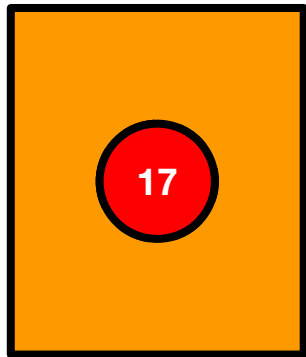
Tree would get
updated
(shown out of
date)



CFS Run Queue
Red-Black Tree

CFS Example

Processor

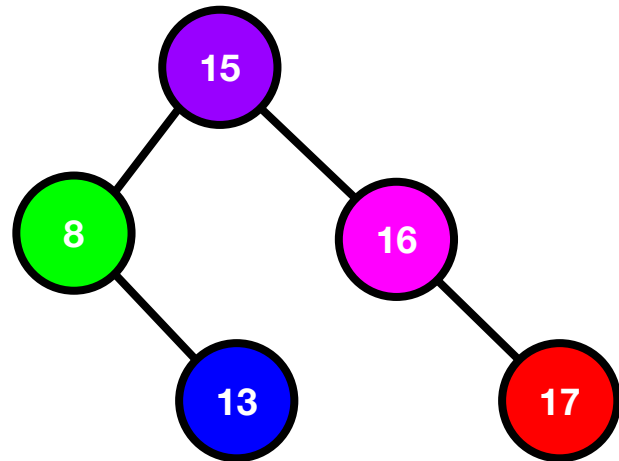


Default nice (nice=0)
Normal priority

Processor Virtual Runtime Timeline



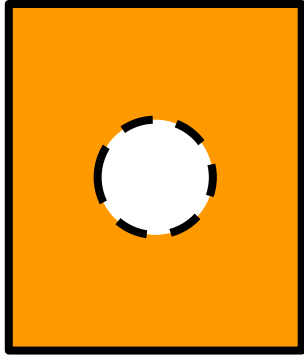
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

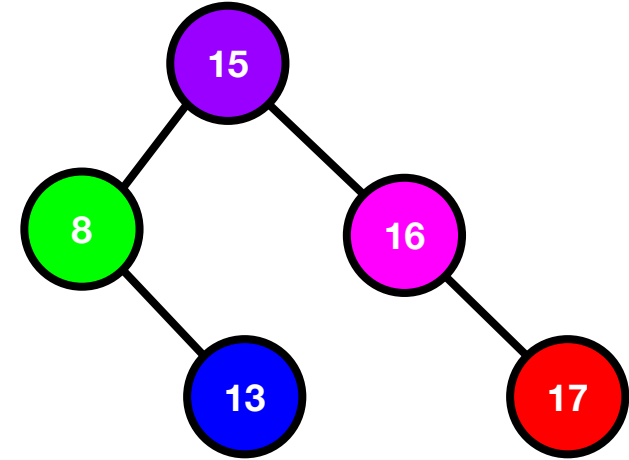
Processor



Processor Virtual Runtime Timeline



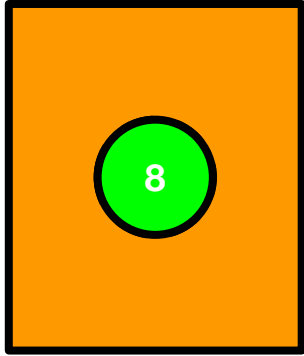
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

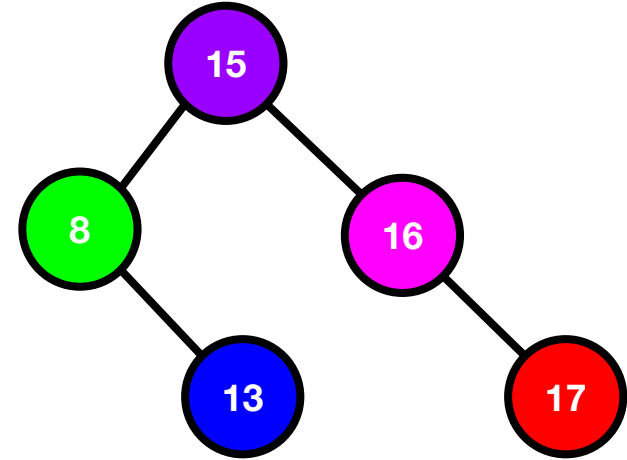
Processor



Processor Virtual Runtime Timeline



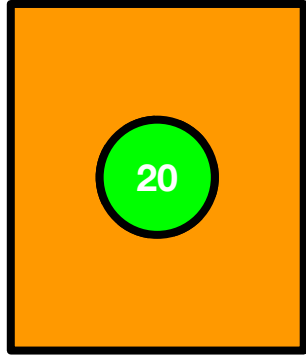
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Processor

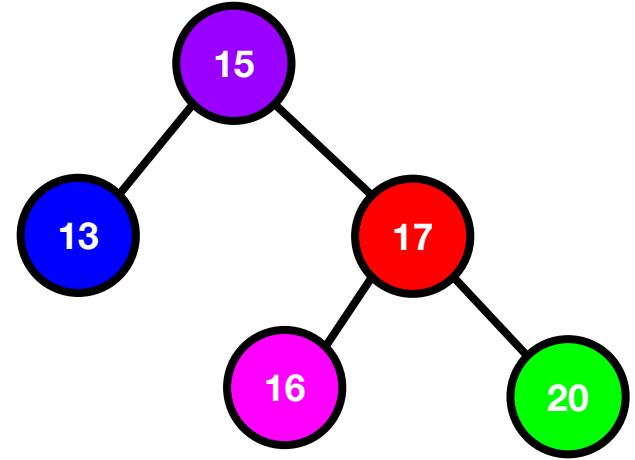


More nice (nice>0)
Lower priority

Processor Virtual Runtime Timeline



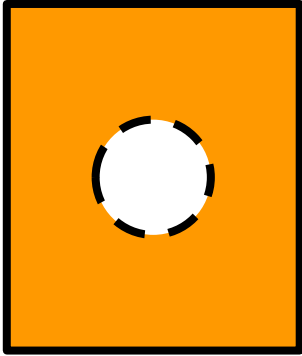
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

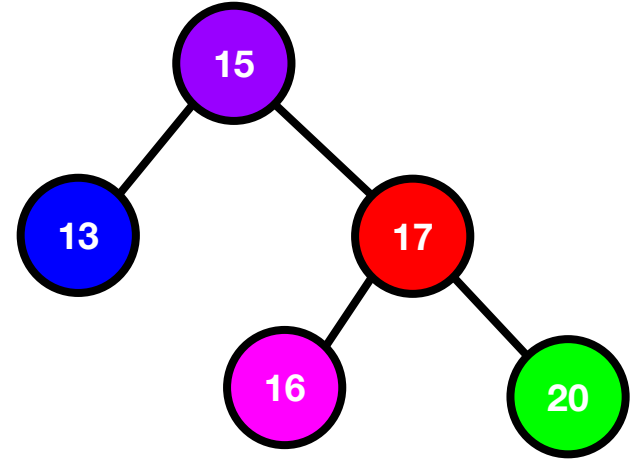
Processor



Processor Virtual Runtime Timeline



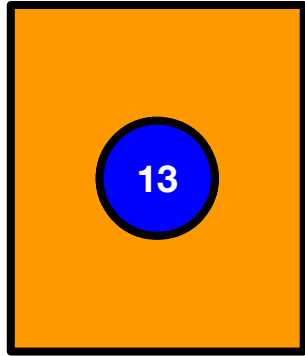
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

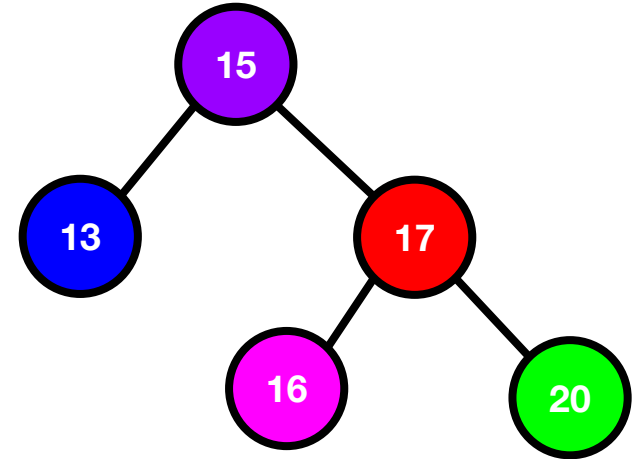
Processor



Processor Virtual Runtime Timeline



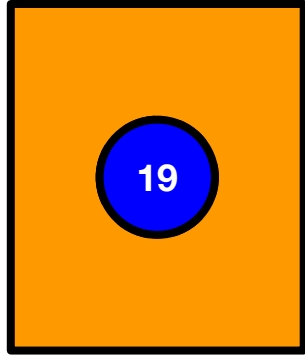
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Processor



More more nice (nice>>0)
Very Low priority

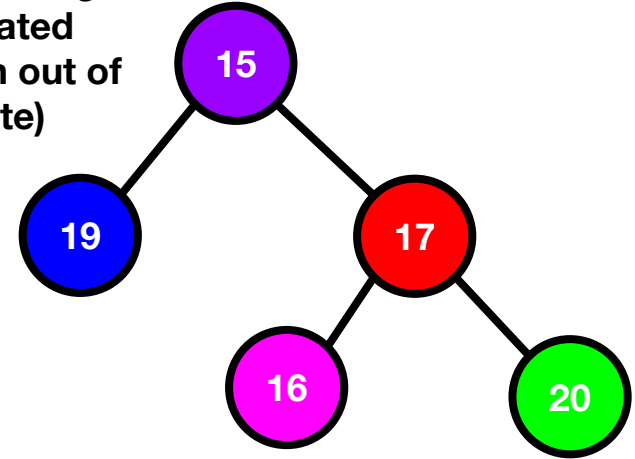
Processor Virtual Runtime Timeline



Processor Wall-Time Timeline



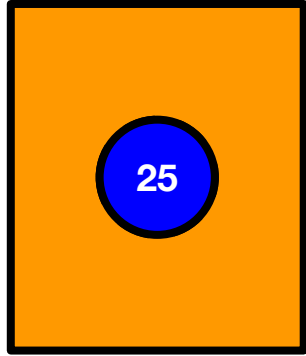
Tree would get
updated
(shown out of
date)



CFS Run Queue
Red-Black Tree

CFS Example

Processor

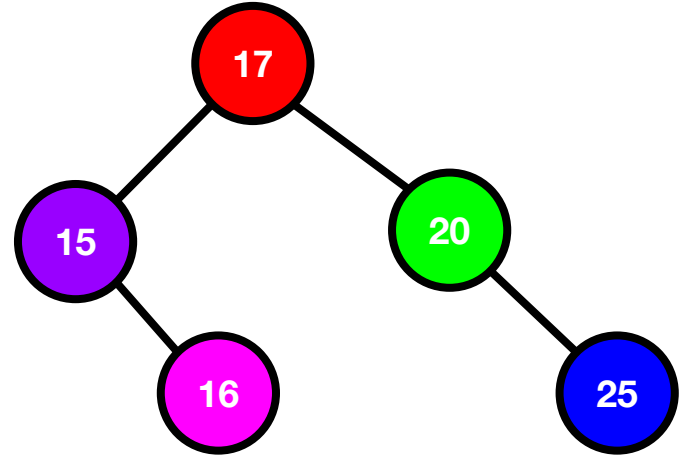


More more nice (nice>>0)
Very Low priority

Processor Virtual Runtime Timeline



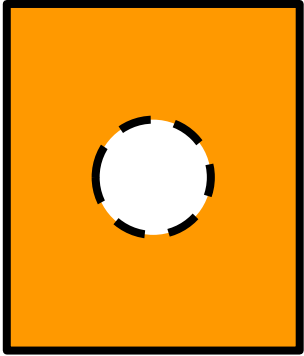
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

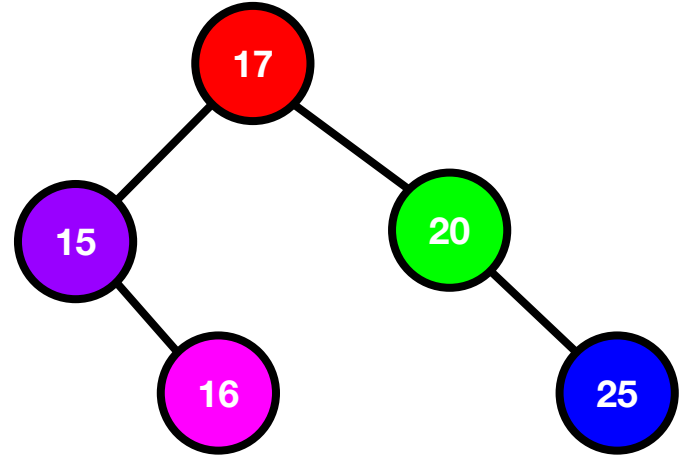
Processor



Processor Virtual Runtime Timeline



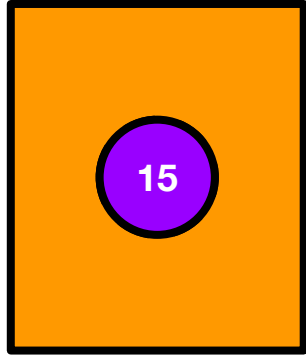
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

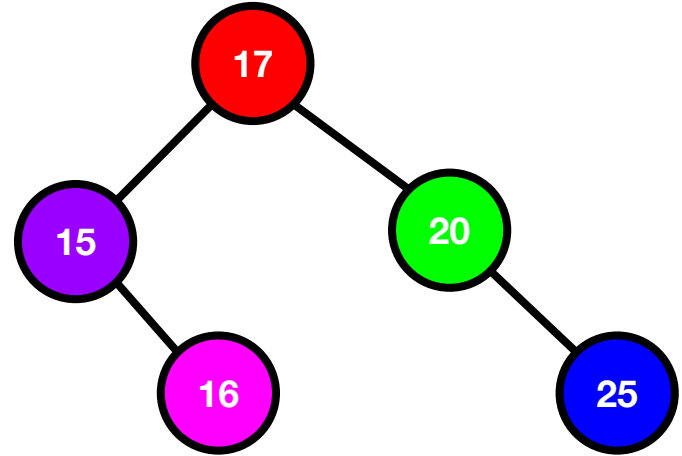
Processor



Processor Virtual Runtime Timeline



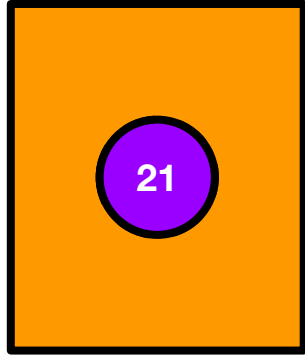
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Processor



Slightly not nice ($\text{nice} < 0$)
Slightly high priority

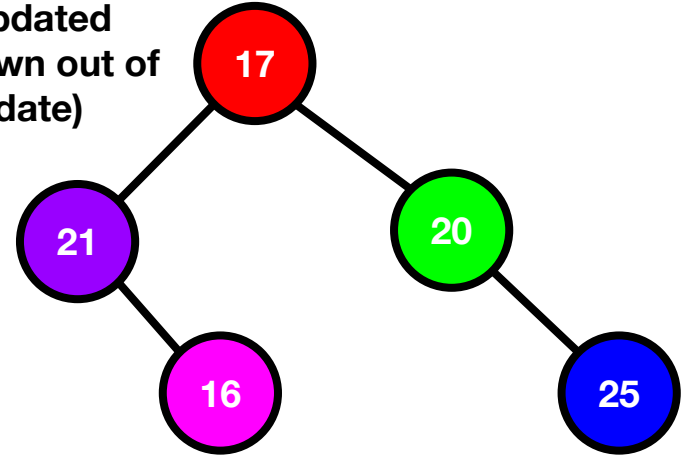
Processor Virtual Runtime Timeline



Processor Wall-Time Timeline



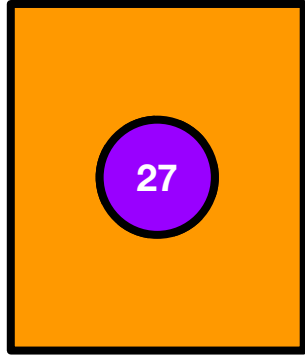
Tree would get
updated
(shown out of
date)



CFS Run Queue
Red-Black Tree

CFS Example

Processor

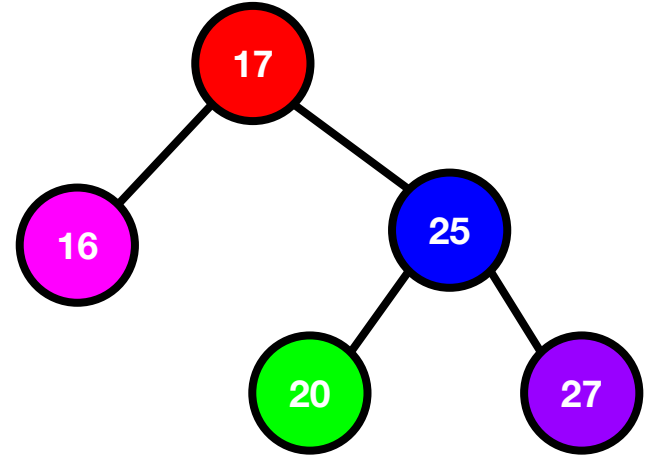


Slightly not nice ($\text{nice} < 0$)
Slightly high priority

Processor Virtual Runtime Timeline



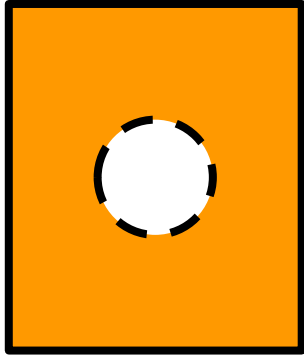
Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

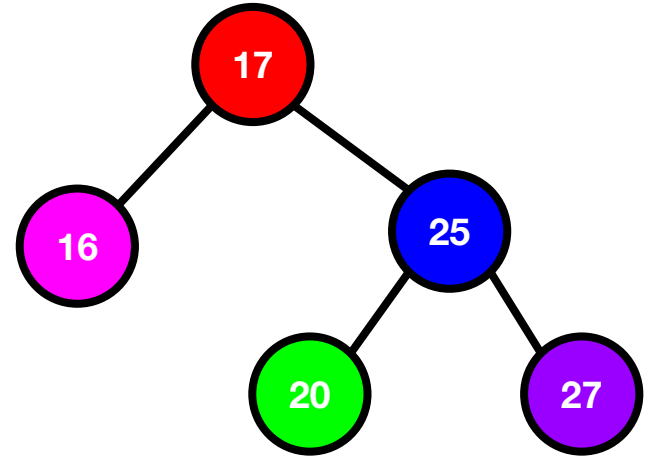
Processor



Processor Virtual Runtime Timeline



Processor Wall-Time Timeline

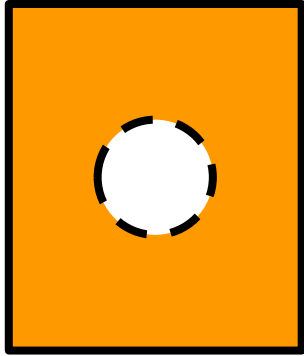


CFS Run Queue
Red-Black Tree

CFS Example

Normalization makes it so all processes' `vruntime` values ultimately have a maximum difference of minimum granularity in virtual time

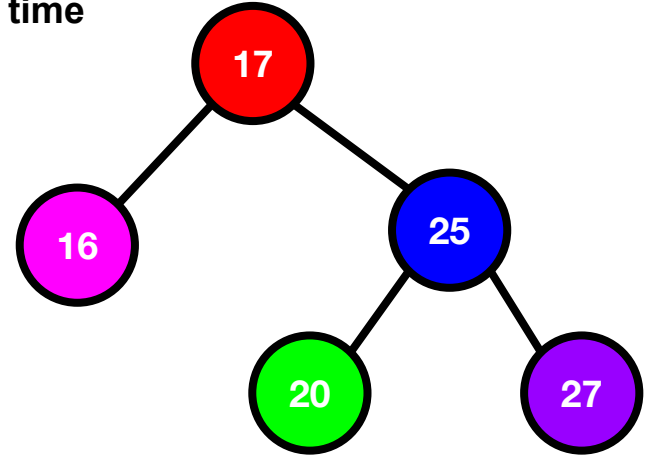
Processor



Processor Virtual Runtime Timeline



Processor Wall-Time Timeline



CFS Run Queue
Red-Black Tree

CFS Example

Normalization makes it so all processes' `vruntime` values ultimately have a maximum difference of minimum granularity in virtual time

Processor

Mike Galbraith

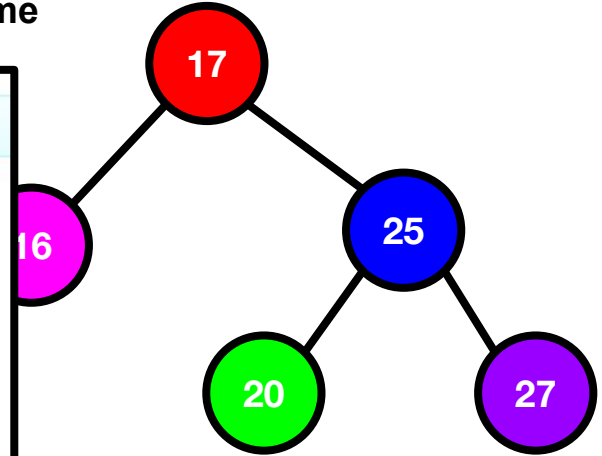
Post by XingChao Wang

Hi Ingo, peter,

When `check_preempt_tick()` selects next leftmost `sched_entity`, it calculates delta `vruntime` of curr and leftmost entity, then compares it with `ideal_runtime`. But `ideal_runtime` is real-time type, need convert it to virtual-time ,right?

Why? The scheduler converges `vruntimes` to within `min_granularity`, that's it's mission. What this test is trying to say is that if the

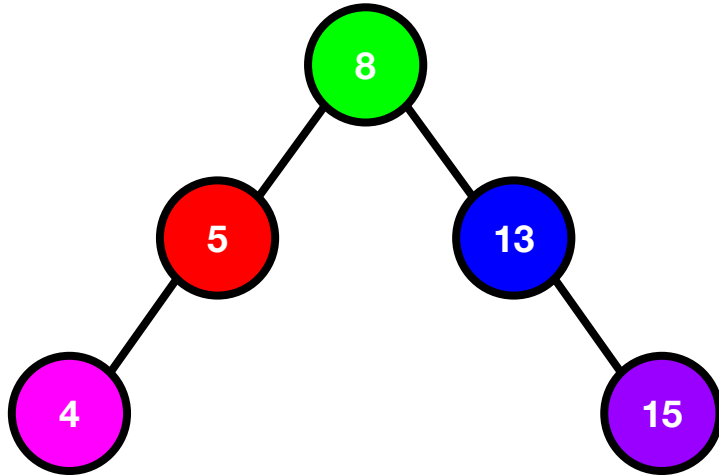
<https://linux.kernel.narkive.com/FKiX036A/check-preempt-tick-check-vruntime-mistake>



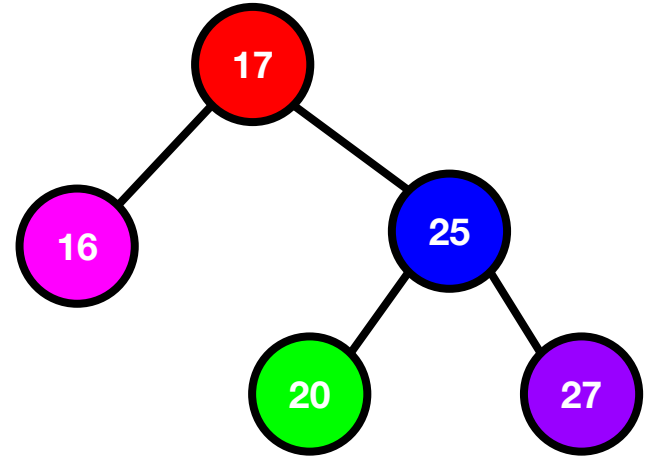
CFS Run Queue
Red-Black Tree

Processor Wall-Time Timeline

Before & After: CFS Run Queue (Red-Black Tree)



CFS Run Queue
Red-Black Tree



CFS Run Queue
Red-Black Tree

How Does The Scheduler Update?

Follow links for info on HR-timer
<https://github.com/torvalds/linux/commit/8f4d37ec073c>
<https://lwn.net/Articles/549754/>

- Periodic timer interrupt
 - 100-1000 Hz (every 1-10ms)
 - In practice, high resolution timers (HR-timer) implemented
- Interrupt calls `scheduler_tick()`
- Statistics about current process computed
- Check if preemption needed
 - Call `schedule()` for deciding next process

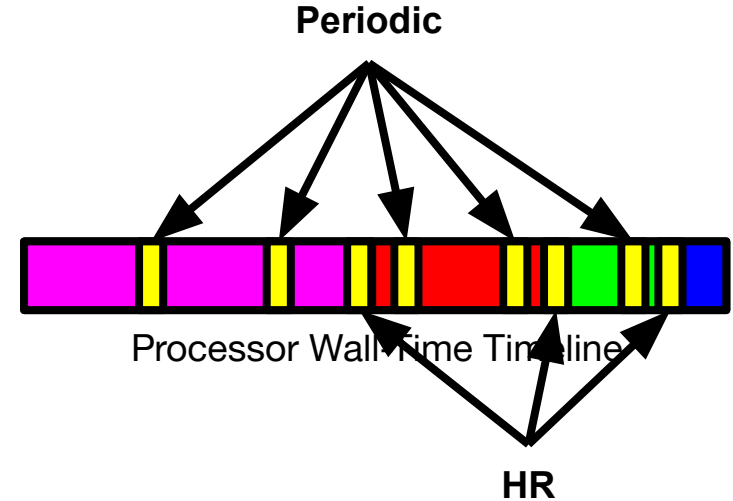


Processor Wall-Time Timeline

How Does The Scheduler Update?

- Periodic timer interrupt
 - 100-1000 Hz (every 1-10ms)
 - In practice, high resolution timers (HR-timer) implemented
- Interrupt calls `scheduler_tick()`
- Statistics about current process computed
- Check if preemption needed
 - Call `schedule()` for deciding next process

Follow links for info on HR-timer
<https://github.com/torvalds/linux/commit/8f4d37ec073c>
<https://lwn.net/Articles/549754/>



Math Behind CFS (3): Increasing Target Latency

- Definitions:

- **sched_nr_latency**: number of runnable processes that should be scheduled in a period
- **sysctl_sched_latency**: scheduling period (ms)
- **sysctl_min_granularity**: supposed wall-time a process would execute

- Consider when the number of runnable processes exceeds `sched_nr_latency`, what happens?

- $\text{period} = (\# \text{ proc}) \times (\text{min_gran})$

$$\text{sysctl min granularity} = \frac{\text{sysctl sched latency}}{\text{sched nr latency}}$$

“Can we get that again in English?”

$$\text{supposed time slice} = \frac{\text{scheduling period}}{\text{number of processes}}$$

$$\text{timeslice} = (\text{scheduling period}) \times \frac{(\text{process weight})}{(\text{sum of all process weights})}$$

“Pop quiz hotshot”



Author: pgrizzaffi, "Pop Quiz, Hotshot. Is It Automated?". From the movie *Speed* (1994). URL: <https://responsibleautomation.wordpress.com/2018/11/01/pop-quiz-hotshot-is-it-automated/>

“Pop quiz hotshot”

Exercise 5

Access Code:



You only need one of these equations - think it over :)

Math time. We have a scheduling period of 20 ms. We have 4 processes: 3 with nice = 0, 1 with nice = 5. What is the wall-time of the nice 5 process? In other words, how long does it *actually* run? How long does each of the individual nice = 0 processes run for?

$$\text{sysctl min granularity} = \frac{\text{sysctl sched latency}}{\text{sched nr latency}}$$

$$\text{timeslice} = (\text{scheduling period}) \times \frac{(\text{process weight})}{(\text{sum of all process weights})}$$

nice	-20	-19	-18	-17	-16	-15	-14	-13	-12	-11
weight	88761	71755	56483	46273	36291	29154	23254	18705	14949	11916
nice	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1
weight	9548	7620	6100	4904	3906	3121	2501	1991	1586	1277
nice	0	1	2	3	4	5	6	7	8	9
weight	1024	820	655	526	423	335	272	215	172	137
nice	10	11	12	13	14	15	16	17	18	19
weight	110	87	70	56	45	36	29	23	18	15

Author: pgrizzaffi, “Pop Quiz, H
<https://responsibleautomation.v>

Nice-to-weight conversion

CFS Implementation - sched_entity

- </kernel/sched/fair.c>
 - Only ~10k lines of code!
- </include/linux/sched.h>
 - struct sched_entity
 - u64 vruntime;
- vruntime - virtual runtime
 - Weighted runtime

```
/ include / linux / sched.h
All sym▼ Search

465 struct sched_entity {
466     /* For load-balancing: */
467     struct load_weight    load;
468     struct rb_node        run_node;
469     struct list_head      group_node;
470     unsigned int          on_rq;
471
472     u64                    exec_start;
473     u64                    sum_exec_runtime;
474     u64                    vruntime;
475     u64                    prev_sum_exec_runtime;
476
477     u64                    nr_migrations;
478 }
```

Linux 5.14.2 Source via Bootlin Elixir Cross Referencer

<https://elixir.bootlin.com/linux/v5.14.2/source/include/linux/sched.h#L465>

CFS Implementation - vruntime & update_curr()

- `update_curr()` keeps vruntime updated
- It also calls `update_min_vruntime()` to keep `min_vruntime` up to date

/ kernel / sched / fair.c

All syn

Search Identifier

```
609 /*
610  * delta /= w
611  */
612 static inline u64 calc_delta_fair(u64 delta, struct sched_entity *se)
613 {
614     if (unlikely(se->load.weight != NICE_0_LOAD))
615         delta = __calc_delta(delta, NICE_0_LOAD, &se->load);
616
617     return delta;
618 }
```

/ kernel / sched / fair.c

All syn

Search Identifier

```
793 static void update_curr(struct cfs_rq *cfs_rq)
794 {
795     struct sched_entity *curr = cfs_rq->curr;
796     u64 now = rq_clock_task(rq_of(cfs_rq));
797     u64 delta_exec;
798
799     if (unlikely(!curr))
800         return;
801
802     delta_exec = now - curr->exec_start;
803     if (unlikely((s64)delta_exec <= 0))
804         return;
805
806     curr->exec_start = now;
807
808     schedstat_set(curr->statistics.exec_max,
809                  max(delta_exec, curr->statistics.exec_max));
810
811     curr->sum_exec_runtime += delta_exec;
812     schedstat_add(cfs_rq->exec_clock, delta_exec);
813
814     curr->vruntime += calc_delta_fair(delta_exec, curr);
815     update_min_vruntime(cfs_rq);
816
817     if (entity_is_task(curr)) {
818         struct task_struct *curtask = task_of(curr);
819
820         trace_sched_stat_runtime(curtask, delta_exec, curr->vruntime);
821         cgroup_account_cputime(curtask, delta_exec);
822         account_group_exec_runtime(curtask, delta_exec);
823     }
824
825     account_cfs_rq_runtime(cfs_rq, delta_exec);
826 }
```

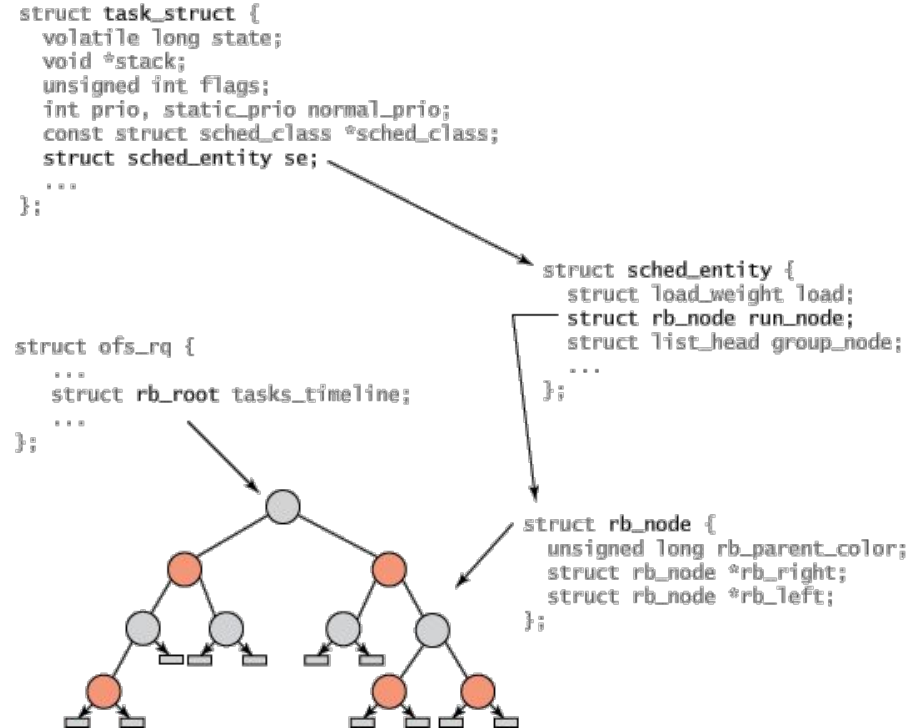
Linux 5.14.2 Source via Bootlin Elixir Cross Referencer

<https://elixir.bootlin.com/linux/latest/source/kernel/sched/fair.c#L793>

65

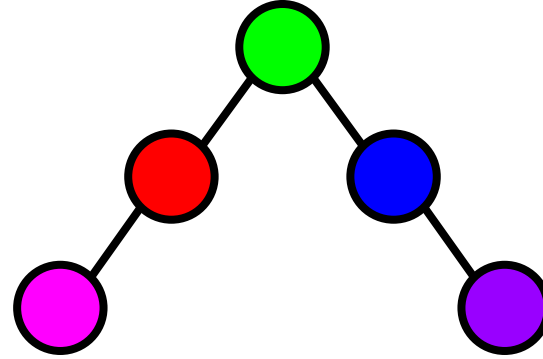
CFS Implementation - Process Selection

- What do we execute next?
- Process having shortest `vruntime`
- This is why we say “fair”
- Red-black tree (rbtree) data structure implemented for processes to sort them by `vruntime`
- rbtree is self-balancing binary search tree

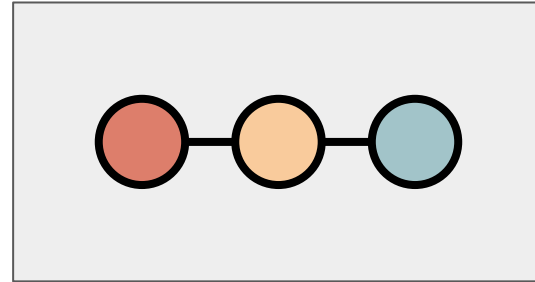


CFS Implementation - Sleeping & Waking Up

- Sleeping (blocked) processes cannot be run
 - Awaiting event
- Events include:
 - Time length
 - File I/O
 - Keyboard input
- Process voluntarily sleeps
 - Moves from run queue to wait queue
 - Invokes `schedule()`
- Reverse this to wake up
 - What will happen when process added to run queue? `vruntime`?



Run Queue ☕
Red-black tree



Wait Queue 🛌
Linked list

Shh! They're **sleeping!**

Context Switching

(Switching Processes)

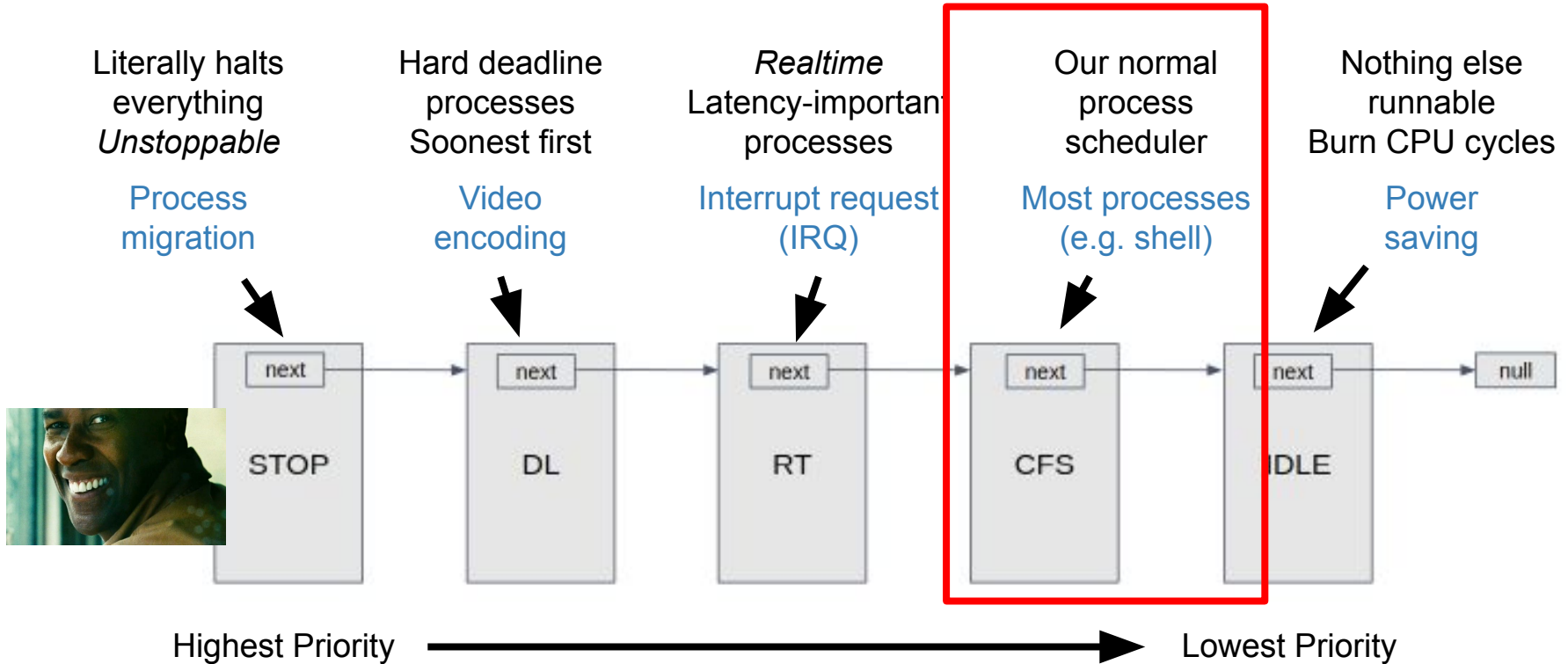
- `context_switch()`
- Switches mm (memory map)
- Switches register states

```
/ kernel / sched / core.c
4640  /*
4641   * context_switch - switch to the new MM and the new thread's register state.
4642   */
4643  static __always_inline struct rq *
4644  context_switch(struct rq *rq, struct task_struct *prev,
4645                struct task_struct *next, struct rq_flags *rf)
4646  {
4647      prepare_task_switch(rq, prev, next);
4648
4649      /*
4650       * For paravirt, this is coupled with an exit in switch_to to
4651       * combine the page table reload and the switch backend into
4652       * one hypercall.
4653       */
4654      arch_start_context_switch(prev);
4655
4656      /*
4657       * kernel -> kernel    lazy + transfer active
4658       * user  -> kernel    lazy + mmgrab() active
4659       *
4660       * kernel -> user     switch + mmdrop() active
4661       * user  -> user     switch
4662       */
4663      if (!next->mm) { // to kernel
4664          enter_lazy_tlb(prev->active_mm, next);
```

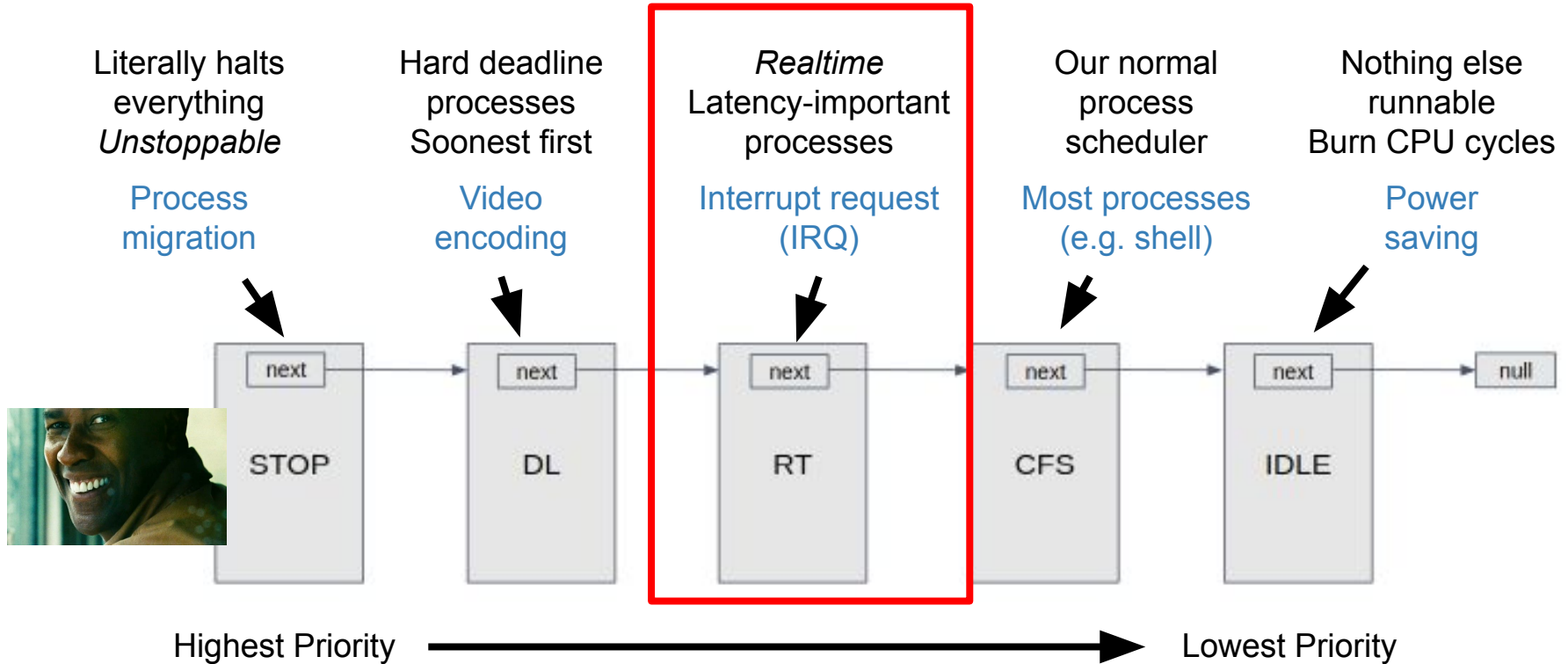
Linux Source via Bootlin Elixir Cross Referencer

<https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L4644>

We Just Talked About CFS (Normal Processes)



Let's Talk About Realtime



Real Time - A Tale of Two Policies

- `SCHED_FIFO`
 - FIFO: First in, first out
 - Process executes until yielding
 - No timeslices
 - Very different from CFS
 - Could block or be done executing
- `SCHED_RR`
 - RR: Round-robin
 - Round-robin with timeslices and priorities

Real Time - A Tale of Two Policies

- `SCHED_FIFO`
 - FIFO: First in, first out
 - Process executes until yielding
 - No timeslices
 - Very different from CFS
 - Could block or be done executing
- `SCHED_RR`
 - RR: Round-robin
 - Round-robin with timeslices and priorities

So we have a beautiful real-time system if we need it?

Real Time - A Tale of Two Policies

- `SCHED_FIFO`
 - FIFO: First in, first out
 - Process executes until yielding
 - No timeslices
 - Very different from CFS
 - Could block or be done executing
- `SCHED_RR`
 - RR: Round-robin
 - Round-robin with timeslices and priorities

So we have a beautiful real-time system if we need it?

Not at all

Real Time - A Tale of Two Policies

- `SCHED_FIFO`
 - FIFO: First in, first out
 - Process executes until yielding
 - No timeslices
 - Very different from CFS
 - Could block or be done executing
- `SCHED_RR`
 - RR: Round-robin
 - Round-robin with timeslices and priorities

The Reality



https://twitter.com/Marjan_Lion/status/727692489264484352/photo/1

Real Time - A Tale of Two Policies

- `SCHED_FIFO`
 - FIFO: First in, first out
 - Process executes until yielding
 - No timeslices
 - Very different from CFS
 - Could block or be done executing
- `SCHED_RR`
 - RR: Round-robin
 - Round-robin with timeslices and priorities

The Reality



https://twitter.com/Marjan_Lion/status/727692489264484352/photo/1

Soft Real-Time

No guarantees on timing

System Calls

Table 4.2. Scheduler-Related System Calls

System Call	Description
<code>nice()</code>	Sets a process's nice value
<code>sched_setscheduler()</code>	Sets a process's scheduling policy
<code>sched_getscheduler()</code>	Gets a process's scheduling policy
<code>sched_setparam()</code>	Sets a process's real-time priority
<code>sched_getparam()</code>	Gets a process's real-time priority
<code>sched_get_priority_max()</code>	Gets the maximum real-time priority
<code>sched_get_priority_min()</code>	Gets the minimum real-time priority
<code>sched_rr_get_interval()</code>	Gets a process's timeslice value
<code>sched_setaffinity()</code>	Sets a process's processor affinity
<code>sched_getaffinity()</code>	Gets a process's processor affinity
<code>sched_yield()</code>	Temporarily yields the processor

Next Lecture:

System Calls