

- How to write a scheduler. Look in few places: (3.10 kernel)
 - ▶ /kernel/sched ← define scheduling class, look at idle/rt/fair.
 - ▶ /include/linux/sched.h ←
 - ① Define a scheduling class.
 - enqueue-task-rt uses sched-rt-entity defined in sched.h also has sched-entity within the task_struct
- ▶ The entity is the thing that is getting scheduled. Rather than schedule tasks directly, you go through entity (which could be used to represent a group of tasks, for example).
An entity eventually goes to a task, since this is what gets ultimately scheduled.
- ② Define an entity for the new scheduler (class specific one) in include/linux/sched.h
 - ▶ Add stuff that your scheduler may need there: ex. timeslice
 - ▶ You will need to link list entities together since you will need to organize them in some structure later (runqueue or sth. like that).
 - ▶ From task-struct, you have pointers to the "encapsulating" entities member vars
- ③ Add the entity member to task-struct in sched.h and make sure the entity can support being in a queue or list (ex. next pointer).
 - ▶ kernel/sched/sched.h ← check out what's there for rt.
 - rt_rq and cfq_rq ← look we have class-related runqs in here!
 - These are the class-specific runqueues.
- ④ Both are part of struct rq ← the main, per-CPU runqueue
- ④ Add a class-specific runqueue for the new scheduling class to the general runqueue and define that new runq. in sched/sched.h
 - ⑤ has a runq lock!
 - THIS IS WHERE YOU KEEP THE ENTITIES FOR THE NEW SCHEDULER.
- NOTE: The RT scheduler defines an array of subqueues "struct rt-prio-array active" in "struct rt_rq" ← one subqueue for each priority level 0..99.

⑤ Add initialization for the new runqueue ~~in rt.c~~ similar to:
init_cfs_rq and init_rt_rq in sched/sched.h.

those are external in sched.h and defined in rt.c, etc.

NOTE: Searching code for "sched-rt" and "-class" (found "sched-class *" in linux/sched.h but said not important)

▶ struct sched_class in sched/sched.h is the general scheduling class. We also have init_sched_rt_class(void)

⑥ Add initialization for the new scheduling class.

▶ Note that in sched/sched.h the ~~classes~~ scheduling classes are listed as "extern const struct sched_class" and there are "sched_class_highest" and "for_each_class" macros

⑦ Add the new scheduling class to the scheduling class list

▶ kernel/sched/core.c is the core area of the scheduler

- Look for class-specific stuff there: "_rt-", etc. see if they matter.

- Look for system calls that set:

- ⊗ scheduling parameters

- ⊗ what scheduling class you are part of

Ⓢ Look for SYSCALL_DEFINE

Ⓢ EX: --setscheduler(...) ←

⑧ Hook up your new scheduling class here, so it gets used

--sched_setscheduler(...) ← gets task, policy, parameters and assigns itself to a scheduling class

⑨ Add new scheduling class here, if necessary.

Ⓢ Basically, find all SYSCALL_DEFINES in sched/core.c and see if they do anything that is scheduling-class-specific and if we need to modify them for the new scheduling class.

► -- sched-fork() - sets up scheduler stuff for new tasks

⑩ Modify sched-fork for new scheduling class, also see sched-fork() w/o -- at start!

NOTE: Here we set the policy constant to SCHED_NORMAL, which is the policy for the CFS class, ~~maybe~~ we need to change that default? (this is done for non-RT tasks) if we want to start w/ the new class.

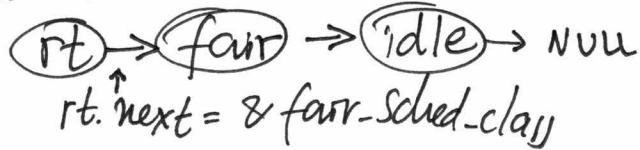
Most of the stuff we want to do is in the class specific file: /kernel/sched/idle-task.c, rt.c, fair.c

⑪ Define a new class (idle-task.c is simplest example but rt.c is more useful)

Some key elements of the class and related functions

→ ④ init-rt-rq - initialize the RT runqueue (the CFS init is much simpler)

→ ⑤ .next is the link to the next scheduling class on list:



Add new class here

b/w rt and fair

⑥ enqueue-task and dequeue-task ← check out RT implementation and CFS implementation

NOTE: "foreach-sched-entity" is for group scheduling, if you don't worry about that this loop gets executed only once.

enqueue-task calls -- enqueue-entity and thus:

enqueue-entity-rt calls → ① gets the rung array
→ ② figures out where to put in the array

Essentially, we find the entity and put it on a rung - the complexity comes because RT allows for groups where entities form a tree-like structure. If we don't worry about this, it is simpler:

- ① take the task, find the scheduling entity associated w/ it
- ② stick it on your runqueue which is a list or sth. like that

dequeue-task is the opposite of enqueue-task

for RT: dequeue-rt-entity → dequeue-rt-stack → dequeue-rt-entity
should be simpler w/o groups

① pick-next-task

for RT: pick-next-task → pick-next-task-rt → pick-next-task-rt

Go thru the runque and find next entity to go to.

② rt-task-of is what is used to find the task of an entity!

HINT: Maybe use a list of entities and don't search the runque
But just pick the entity at the head and then get its task and return it.

③ You probably want to implement everything that is implemented in idle-task.c - don't need to actually do work, but set the function pointer. Ex: prio-changed and switch-to do nothing for idle, but the functions are defined!

NOTE: prio-changed has code for fair.c but if you don't have priorities, then you probably don't have to worry about it (except define & set pointer)
switch-to has code for fair.c - maybe could be left empty or the reschedule flag could be set
get-r-interval - nothing for idle.

④ → Do we need to do nothing/return 0 here or something else?

⑤ check-preempt-curr - for idle calls resched, since idle tasks are unconditionally rescheduled
Do we do nothing?

put-prev-task - called right before pick-next-task in core.c.

⑥ Do we need to do sth. to the task that was already running before picking the next task? **PROBABLY!**

set-curr-task - nothing for idle, for fair.c it "accounts for a task changing its policy or group." Gets called when a task migrates b/w classes.
If a task is currently running & its scheduling parameters change (ex. scheduling

class changes, do you need to do anything special for it?

PROBABLY!

◀ dequeue, enqueue, and pick-next-task are the functions that must be defined for sure!

ⓕ task-tick is the other thing you probably need.

Ⓢ gets called on a scheduler tick. task-tick-rt finds the scheduling entity, then (for RR tasks only, FIFO have no slice) decrements the time slice, checks to see if it has expired and, if so, update the timeslice back to starting count, and re-enqueue to end of runq and set-tsk-need-resched() since task-tick gets called in the context of a timer interrupt, so we can't call schedule() directly.

ⓖ yield-task - gets called when a task yields the CPU: for RT what we do is requeue the task

ⓗ SMP stuff (ifdef CONFIG-SMP)

Ⓢ select-task-rq → nothing for idle (idle tasks not moved)
 Ⓢ pre-schedule → for RT: gets called from fork, wakeup, i.e. when a task becomes runnable and we need to pick a CPU to put the task on its runq.
 Ⓢ post-schedule → LOCKING is required & think about dead/lock!

Ⓢ Do we need this?

Ⓢ Called before pick-next-task, toward the beginning of schedule.
 Ⓢ Do we need to do sth. before pick-next-task that is SMP related?

ⓗ HINT: More things between runqueues, for example, before picking the next task to run!

SUMMARY:

- ① linux/sched.h - add entity stuff (definition & task-struct member)
- ② sched/sched.h - define custom runq & insert in general runq
- ③ core.c - check system calls and the underlying set_scheduler & set_sched_scheduler that set scheduling parameters & set scheduling policies
- ④ class-specific.c - implementation of scheduling class:
 everything in idle-task.c + enqueue + yield-task + task-tick + select-task-rq + pre-schedule + post-schedule (?)

- Group scheduling - The "for-each-..." macro in scheduler code has to do w/ group scheduling
- An entity may have several tasks - i.e. a group of tasks are scheduled as a unit. Ex: `while(1) fork(2);` will kill system if scheduled as separate entities w/ separate time quanta. With group scheduling, we have the original sched_entity and all subsequent tasks will be part of it (as a tree). When we schedule we only look at that one group.
- The "for-each..." gets to the leaves in the entity group tree and gets their tasks.
- sched_entity has a parent if we use groups, if parent is NULL, then the "for-each..." loop is executed only once.

● MEMORY: ▶ processes have virtual addresses that refer to data:



Benefits: Conserve space: say you do `malloc(20GB)` and `exit()` - the OS would not really reserve 20GB of physical memory for you.

▶ When and How does mapping occur?

- Compile-time ← you would map all virtual addresses to physical ones, but that would not work well w/ multiple processes, also, this is not feasible for dynamic allocation.
- Load-time ← see what is free at time of loading and then map to physical addresses that are free. Still not feasible for dynamic allocation & not efficient.
- Execution-time ← as the program is running, the mapping b/n virtual & physical memory is adjusted, if necessary. Most flexible but overhead of translation for each memory access is not feasible unless hardware support is available.

Needs hardware support!
(ex. paging)