

Article SaaKM

Dubois Swan Gouicem Redha Pires Baptiste Sopena Julien

August 12, 2025

1 Abstract

2 Introduction

CPU scheduling is a fundamental aspect of operating systems and plays a critical part regarding performances and security. With the increasing complexity of hardware and the diversity of workloads servers have to run, writing schedulers is tedious task.

3 Background

The scheduler is a critical part of an operating system (OS) that manages the processes executions. It is required to take into account the hardware it runs on and it might need to adapt to a lot of different workloads. It is also tightly coupled with other parts of the OS like the memory management subsystem (especially regarding NUMA architectures), it therefore require a lot of knowledge to write and maintain a scheduler. When it comes to the Linux kernel, the Earliest Eligible Virtual Deadline First (EEVDF) scheduler is the default general purpose scheduler since v6.6. It is composed of 7000+ LoC alone and depends on a significant amount of C headers and kernel APIs. Because of its genericity it implements a lot of heuristics that can be hard to isolate and understand. It also causes a loss in performances in order to adapt to various workloads. That is why there is an effort in the industry and academia to offer kernel developers a way to write and test schedulers easily *[ghost, schedext, plugsched]*. This need motivated and allowed the successful merge in the kernel of *sched_ext* in the Linux version v6.12.

In this section, we will go through the current options that are available to kernel developer to write and test schedulers in the linux kernel. First, we will see the *sched_class* API of the Linux kernel, then we will take a look at Google proposition, ghOst[ghost-paper] that delegate

scheduling decisions to user space. Then, we will see *plugsched* a solution proposed by Alibaba, finally we will see *sched_ext* the latest solution proposed by Meta that has been merged into the Linux kernel.

The *sched_class* API *[tableau comme CSD avec +/- de chaque solutions ?]*

Linux scheduler subsystem offers an API to write scheduler through its *sched_class* structure. It contains a set of function pointers that has to be implemented by the schedulers and we will be called in the right place by the scheduler core. You still need to modify some datastructures of the kernel like the *task_struct* (it represents a task) and the *rq* (it represents a runqueue) to store data of your scheduler. This already represents a lot of work because you need to understand how the core of the scheduler (6000+ LoC) works, when and how to take locks and all of the other synchronization mechanisms used (RCU, memory synchronizations, ...).

Once you have implemented your scheduler, you need to modify the linking file *vmlinux.lds.h* to add your brand new scheduling class. This is what will allow the kernel to find it and register it with the others. The last step is to compile and boot your kernel, and only then, you will be able to test and debug it.

In short, writing a scheduler using the *sched_class* API is a long and tedious process that requires a wide and deep knowledge of the Linux kernel. *[réécrire + restructurer un peu cette partie avec schedma sched class + avantages et inconvénients]*

ghOst - Userspace scheduling In order to match the evolution of their needs, Google needed a way to write, test and deploy schedulers easily. That is the main motivation behind ghOst[ghost-paper]. It delegates scheduling decisions to userspace and allows to deploy new schedulers or features without having to recompile and deploy kernel images. ghOst architecture is divided into two parts. The first one lies

in the kernel, using a *sched_class* to implement the core of ghOSt and to define a rich API exposed to the userspace. The second part is executed in userspace and is composed of *agents* which are scheduling policies. The two parts communicate through *message queues* that allows the kernel to notify the agents of a thread state change, a scheduling tick. The agents now have all the information needed to make scheduling decisions and can commit *transactions* via shared memory. [\[presentation resultats + limitations\]](#)

plugsched - Live scheduler update

Plugsched takes another approach to give users a way to tweak and modify the scheduler. It leverage the modularization of the Linux kernel by isolating the scheduler code into a *Linux kernel module* that users can then modify, recompile and insert it into a running kernel. To do so, they rewrite the code sections of threads to redirect modified functions. [\[on est forcés de de passer tous les threads dans la nouvelle politique\]](#)

sched_ext *sched_ext* is a proposition by Meta that has been merged into the Linux kernel mainline in v6.12. It offers a way to write schedulers as *eBPF* programs.

As we have seen in this section, the need for a way to write, test and easily deploy new schedulers in the Linux kernel is a pressing issue. Each of the presented solutions have their own advantages and drawbacks

4 Scheduler as a Kernel Module

We now present *Scheduler as a Kernel Module* (SaaKM), a framework that allows to write schedulers as Linux kernel modules. Our main goal is to provide kernel developers a way to write, test and deploy schedulers easily. To do so, we hide the complexity of the core scheduler (synchronization mechanisms, ...) behind a set of functions corresponding to scheduling events that each policy must implement. We are also capable to have multiple policies loaded at the same time, allowing each applications to chose the scheduler that best fits its needs.

Design We rely on the *sched_class* API presented above to implement a minimal core scheduler (less than 1500 LoC) that will only be used to register policies and call the callbacks at the right places.

Each policy must implement a set of functions that map to specific scheduling events (e.g wakeup, scheduling tick, ...) and then register itself to our scheduling class.

Events The events are divided into two categories :

- Task events are all events related to task, such as : a blocking task, a task waking up.
- Core events are all events related to the CPUs, such as : a scheduling tick, a CPU becoming idle.

In order to hide the complexity of the locks mechanisms used in the core scheduler, we deviated