

SaaKM : Schedulers as Kernel Module*

*Note: Sub-titles are not captured in Xplore and should not be used

1st Baptiste Pires
LIP6
Sorbonne Université
Paris France
baptiste.pires@lip6.fr

2nd Given Name Surname
dept. name of organization (of Aff.)
name of organization (of Aff.)
City, Country
email address

3rd Given Name Surname
dept. name of organization (of Aff.)
name of organization (of Aff.)
City, Country
email address

4th Given Name Surname
dept. name of organization (of Aff.)
name of organization (of Aff.)
City, Country
email address

Abstract—*todo:[abstract]*

Index Terms—kernel, scheduler, module

I. 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. Linux provide a general purpose scheduler, Earliest Eligible Virtual Deadline First (EEVDF) that is designed to handle those constraints. However, due to its genericity, it is not tailored for specific workloads and offers worse performance than tailored scheduling algorithm. If you want to tweak it to better fit your workloads, you need to have a deep understanding of the core scheduler. Furthermore, the process to update, compile, deploy and test the new version is time consuming. That is why there are efforts to provide kernel developers ways to write, test and deploy schedulers faster.

II. MOTIVATION AND BACKGROUND

The scheduler is a critical part of an Operating System (OS) that manages processes execution. It has to take into account the hardware it runs on but also have to adapt to a lot of different workloads. Currently, the general purpose scheduler of Linux is Earliest Eligible Virtual Deadline First (EEVDF) that has been introduced in v6.6. Being the default scheduler, it has to support a wide variety of workloads (e.g video encoding, web services, ...) but this genericity comes with a cost. Firstly, it is composed of 7000+ LOC and depends heavily on other kernel APIs like the memory management subsystem, the synchronization mechanisms (e.g locks, RCU, ...). This make the development and maintenance quite complex for the kernel developers as you need to put

a lot of time and effort to be able to understand the code. Being a core part of the kernel, modifying it requires you to recompile and redeploy your kernel, making the development and debugging a time consuming process. That's why there are been efforts to provide kernel developers ways to write, test, debug and deploy schedulers easily.

There has been an effort in recent years to provide such a solution. We can divide those in two main categories, the schedulers executing in userspace and the ones executing in kernel space.

A. Userspace schedulers

Userspace schedulers are executed outside of the kernel. They require a library to be linked against that will handle the communication with the kernel.

ghOSt - Userspace agents Google proposed ghOSt [?] is an infrastructure that delegates scheduling decisions to userspace. It is composed of two parts. The first one lies in the kernel and implement a scheduling class that interfaces itself with the core scheduler. The second part is executed in userspace and is made of agents. Agents implement scheduling policies thanks to a rich API exposed by the kernel part. The scheduling class exposes scheduling information through messages queues then they inform the kernel part of ghOSt of their scheduling decisions by committing transactions via shared memory.

Skyloft - Userspace interrupts (impossible de trouver le code alors qu'il y a des artifacts ?) Skyloft takes yet another approach at userspace scheduling. They isolate CPUs core to be exclusively used by their userspace scheduler expect for some background kernel threads. They create kernel threads per isolated cores that are bound to these and a kernel

Identify applicable funding agency here. If none, delete this.

module exposes an ioctl interface for userspace program to communicate its scheduling decisions. They achieve userspace us-scale scheduling preemption by leveraging usqerspace interrupts (UINT), thus bypassing the kernel. It allows interrupts to be delivered to userspace without swapping adress space or privilege level.

It also exists library that allows to do userspace scheduling without having to modify the kernel such as ... but they are out of the scope of this paper as they rely ...

B. Kernel schedulers

Thread scheudling being one of the most important part of an kernel, most efforts forward the ease of scheduler development focus on the kernel space.

Scheduling class 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 need to modify some scheduler related datastructures to store data of your scheduler. This already represents a lot of work and time because you need to understand how the core of the scheduler (more than 6000 LOC) works. It heavily relies on complex synchronization mecanismes such as locks, Read-Copy-Update (RCU) and barriers due to the inherent concurrent nature of a multicore architecture. Once the scheduling class is implemented, you need to modifu the *vmlinux.lds.h* linking file to add your scheduling class. This allows the kernel to load and register it in its list, as shown in Figure ?? showing all of the current scheduling classes in the Linux knerel. You then need to modify the scheduler init function to modify some checks that ensures that all scheduling classes are loaded properly. Finally, you can compile and boot your kernel. Only then, you will be able to test and debug it. Writing a scheduler with the scheduling class API is a long and tedious process that requires a wide and deep knowledge of the Linux kernel that can be a barrier for developers.



Fig. 1. Linux scheduling classes overview

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. Once you have modified the extracted code, you can insert it. They do a data update during the insertion to migrate the data from previous scheduler version to the new one.**note:[reformuler, on est forcés de de passer tous les threads dans la nouvelle politique]**

Enoki - Enoki exposes

sched_ext - eBPF schedulers The ext scheduling class is a proposition by Meta that has been merged into the Linux kernel mainline in version v6.12. It offers a way to write schedulers as *eBPF* programs. Their goal is to be capable of quickly test and tweak schedulers. The eBPF verifier provides securtiy and safety as it ensures that the program is correct before loading it. It can also be an issue because of the limitations it imposes, policies must heavily rely on helper functions that they defin**note:[retrouver la ref dans enoki je crois?]**. The eBPF API allows to communicate with userspace program through maps, making it possible to delegate scheduling decisions to userspace. For instance, in their scheduler *scx_rusty*, where they offload all of the load balancing logic to a userspace program written in Rust.

There are others way to modify the scheduler in the kernel. kpatch [] provide a way to live patch a running kernel but is restricted to function level.

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 summarized in Table ?.

III. DESIGN

In this section we present the design of SaaKM, our framework for writing schedulers as Linux kernel modules. First we present the goals of our approach, then we detail the design architetur.

A. Design goals

Writing schedulers in Linux is a tedious task so one of our main goals is to provide a simpler way to do that. The need to have a deep and wide knowledge of the Linux kernel can be a barrier also so we want to hide as much complexity as possible. Finally, testing and debugging a scheduler is also a criteria.

Ease of Use An ideal solution would not require kernel developers to master a new langage nor a new subsystem of the kernel as it would defeat one the purpose of the solution, that is the accessibility part. By the nature of kernel development, they must already be knowledgeable in C and some of the kernel APIs and subsystems. Linux kernel modules answer perfectly to this need.

Performance The solution must incur the lowest overhead possible. As the scheduler is a critiical part regarding performances, even a small overhead can cause significant performances degradation. We must keep our solution as lightheiwt as possible, giving the users the choice to add features that can cause this overhead if they want to.

Testing and Debugging The current process of debugging and testing a scheduler class require to recompile and redeploy

your kernel image each time you make a modification. This can be quite time consuming and slow down the development process. That is why our solution must not be quick and easy to test without the need reboot your machine.

B. Design overview

SaaKM is an API to write schedulers as Linux kernel modules. It exposes a minimal set of exported functions to the modules that can be used to allocate, manage and destroy scheduling policies and data. Figure ?? shows how SaaKM integrates itself with the current scheduling classes. We position ourselves right before the idle class to not disrupt higher priority classes that other threads runs.

All of the SaaKM design relies on a structure composed of handlers that each policy must implement. Each handler maps to a scheduling event. Events are divided into two categories. Thread events correspond to all scheduling events related to a thread (e.g thread is created, waking up, ...). On the other hand, core events are related to the CPU cores (e.g scheduling tick, core becoming idle, ...). This distinction makes it clear for the user to know what it should do, hiding the complexity of the scheduling class API where a single function can be called from multiple paths (i.e enqueue_task). Table ?? shows an excerpt of those events.

Furthermore, this division allows us to hide the complexity of the synchronization mechanisms of the core scheduler. As multicore and NUMA architectures are the quite common now, the scheduler must synchronize data across cores quite often as two core may access the same data. It protects these data through the usage of locks, RCU and low level synchronization mechanisms like memory barriers. Thanks to our design, the user does not need to worry about these and assume that it has the right locks on the data whenever it is needed. **note:***[réécrire cette partie]*

To increase the flexibility of our solution, we support to have multiple policies loaded at the same time. This allows users to select the scheduling policy that is tailored for their workload. **note:***[completer]*

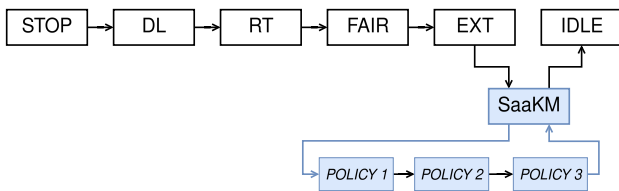


Fig. 2. SaaKM scheduling class with three policies loaded (refaire pour rendre plus visible)

C. Implementation

We implement SaaKM on Linux v6.12, the latest longterm release. In this section we will go through the implementation details.

We implement a minimal (less than 1500 LoC) scheduling class that implements all of the scheduling functions required by the core scheduler. We place ourselves between the ext and idle scheduling classes as shown in Figure 2 where there are three policies loaded. We define a structure `saaKM_module_routines` that is composed of all the handler a policy must implement. Table ?? shows an excerpt of these handlers.

Each policy must allocate and populate a structure with its own handlers and they will be called by our scheduling class at the right places. Each handlers has at least one argument that is a pointer to its corresponding policy structure. Then depending on the event type of the handler (thread or core), it will have either a pointer to a structure `process_event` or `core_event` that contains the data related to the event.

Runqueues management We provide a minimal runqueue definition that policies uses and includes in their own runqueues as a structure field. This allows to delegate the insertion/deletion and mandatory runqueue metadata fields to be handled by SaaKM. We provide three runqueue types. A FIFO, LIST and RED-BLACK tree. You must provide an ordering function for each type of runqueue, except for the FIFO that will always insert the task at the end of the queue. We provide these three types as they are the most common ones but thanks to its design it can be extended. Policies are required to initialize their runqueues before registering themselves.

Thread states We define our own thread states to make to not interfere with the core scheduler. We define 7 states summarized in Figure ?. We have one special state, `SAAKM_READY_TICK` that allows policy to notify SaaKM that it needs to trigger a reschedule. With this states, we are able to check for wrongful transitions and provide debug feedback to the user at runtime. This ease the debugging and testing process. **note:***[peut etre faire un diagramme d'etat qui montre la machine à etat et que c'est ça qui nous permet de detecter les erreurs ? mieux que tableau ?]*

Thread Migration To handle thread migration that occurs because of affinity change or a `exec` rebalance (a newly created thread is being run on a different CPU than its parent), we introduce a new flag, `OUSTED`. This flag allows us to know that we are being dequeued because of those reasons so we can simulate a fast block/unblock sequence to remove the thread from its current runqueue and insert it into the new one.

Load balancing A central and critical part of scheduling is the load balancing. It consists of trying to distribute the load across CPU cores according to some metrics. To do so, you may need to access data belonging to other cores. Locks are therefore needed to ensure the correctness of the data. We identified two types of load balancing. Periodic load balancing occurs at a fixed time interval and idle balancing. We define a new software interrupt, `SCHED_SOFTIRQ_IPANEMA` with

a handler that calls policies `balancing_select` function. Policies can then chose at which time interval they want to do load balancing. We also raise our interrupt when the fair scheduling kicks a CPU to do the NOHZ balancing.**note:[est ce qu'on doit garder cette feature avec EXT ?]**

Policies should also be able to do load balacing when a CPU is becoming idle. Whenever that happens, we call the policy handler to let it perform it.

Hardware topology To allow policies to properly make load balancing, we export a per-cpu variable, `topology_levels` that represent the underlying hardware topology. To do so, we rely on the already existing scheduling domains that are build at initialization of the scheduler subsystem. Policy can then take educated decisions based on the underlying hardware.

Multi-policy support SaaKM supports to have multiple policies loaded at the same time. When it registers, it is given an unique incremental id that is used by application to select their policy. We then maintain them in a linked list ordered by insertion order. When a new thread must be elected, we go through the list of policies starting from the head of the list. The first policy to return a non NULL thread stops the search and the thread is elected. For load balancing, we also go through the full list to make sure that all policies can perform their load balance.

TABLE I
SAAKM THREAD STATES

| State | Meaning |
|------------------|---------------------------------------|
| SAAKM_NOT_QUEUED | Not in a runqueue |
| SAAKM_MIGRATING | Being migrated |
| SAAKM_RUNNING | Running on a CPU |
| SAAKM_READY | Ready to run and in a runqueue |
| SAAKM_READY_TICK | Became ready from a <i>tick</i> event |
| SAAKM_BLOCKED | Blocked and cannot run |
| SAAKM_TERMINATED | Dead |

D. SaaKM workflow

In this section we will go through SaaKM workflow

Registering a policy To register a policy, during its initialization, a module must allocate and fill a `saakm_policy` that contains all of the policy metadata (i.e. name, routines, id and metadata). It must also allocate then initialize its runqueues with the `saakm_runqueue_init` specifying the type of the runqueue (FIFO, LIST, RED-BLACK TREE) and an ordering function if needed. It can then register itself by calling the `saakm_add_policy`.

Selecting a policy Selecting a policy is pretty straighforward. An application must define an `sched_attr` structure and set the `sched_policy` field to `SCHED_SAAKM` and set the `sched_saaKM_policy` with the value of the policy id it wants. From this point, the

thread `task_struct.sched_class` field is set to our scheduling policy and will therefore be scheduled by SaaKM.

Debugging and testing Being a kernel module, policies can leverage the already existing kernel debugging tools such KGDB, ftrace, the debugfs interface and many others. On top of that, we provide a kernel configuration options to check the transitions of the thread states and trigger panic when it happens. We also exposes a `procfs` directory that list the current policies loaded with some metadata. Finally, we use the `sysfs` interface to turn on or off some debugging feature : thread transitions checks, print the wrongful transitions and print the scheduling classes calls.

TABLE II
EXTRACT OF SAAKM_MODULE_ROUTINES FUNCTIONS

| Function | Description |
|------------------------------------|--|
| Thread events | |
| <code>new_prepare(p)</code> | Return CPU id where p should run |
| <code>new_place(task, core)</code> | Insert p into core runqueue |
| Core events | |
| <code>schedule(core)</code> | Called when a task must be elected |
| <code>newly_idle(core)</code> | Called right after <code>schedule</code> |
| Policy management events | |
| <code>init()</code> | Called to initialize policy |

IV. EVALUATION

In this section we present our experimental setup

A. Experimental Setup

We run our benchmarks on a server running Debian 12 bookworm with a our patched Linux kernel v6.12. The server is equiped with 10 core and 20 SMT threads (Intel(R) Core(TM) i9-10900K CPU @ 3.70GHz), 20MB of L3 cache and 64GO of RAM. We run each benchmark at least 50 times and present the mean and standard deviation.

We compare ourselves to the ext scheduler. We implement a minimal FIFO scheduler for comparaisn. Functions to select CPUs are based on threads PIDs as e want to mesaure the overhead of SaaKM and not the efficiency of a scheduling algorithm.

We ran a total of X applications, using the phoronix test suite and builtin benchmarks from applications.

B. Evaluation goal

The goal of our experimentation is to evaluate if SaaKM is a viable solution for writing schedulers and how it performs compared to existing solution. We are motivated by the

- Définir les métriques
- Définir ce à quoi on se compare
- Présenter l'env et les benchmarks et les motiver
- Présentation et analyse des résultats

V. FUTURE WORKS

Scheduling de policy / ajout de prio entre policy ?

Rajouter une API pour faciliter la communication avec userspaces vu que c'est en vogue ?

VI. CONCLUSION

Through this work we show that