# eBPF, plugin and Rust

# 1. Userspace Rust

By userspace Rust, Rust code that runs as a wrapper of the eBPF. For this part any eBPF C code can be used as the kernel part, the lottery scheduler is chosen here.

```
build.rs              -> build script
Cargo.toml
/src
-- bpf_intf.rs        -> shared by C headers
-- bpf_skel.rs        -> skeleton bindings
-- main.rs            -> userspace program
-- stats.rs           -> stats collector
-- /bpf
   -- scx_base.bpf.c  -> eBPF C program
   -- scx_base.h      -> eBPF C header
```

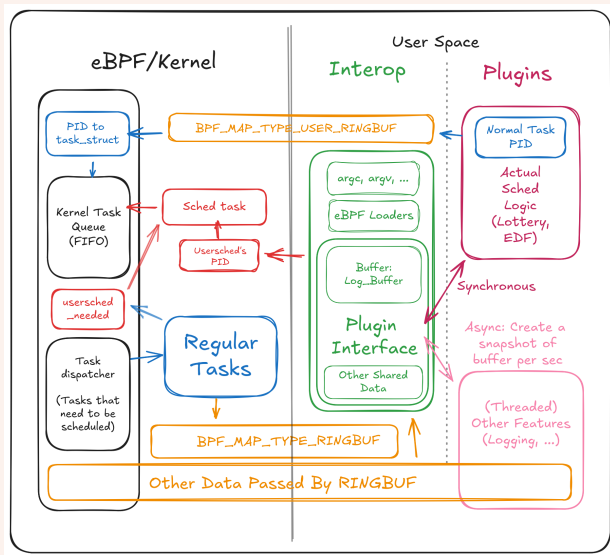Rust uses an entirely different build system to build the eBPF.

```
scx_cargo::BpfBuilder::new()
    .unwrap()
    .enable_intf("src/bpf/scx_base.h", "bpf_intf.rs")
    .enable_skel("src/bpf/scx_base.bpf.c", "bpf")
    .build()
    .unwrap();
```

Rust here serves simply as a loader of the eBPF program, and a collector of statistics, and it does have some advantages compared to C.

- Better ecosystem:
    - Use existing packages for argument parsing, data structures, etc.
    - e.g. Using `circular buffer`, `serde` and `csv` for more flexible logging.

- Better safety and error handling, better atomic variables and strict memory orders all built in, e.g. `Arc{AtomicBool}`, `Ordering::SeqCst`.

- Better build system, without the conflict of header files, especially preprocessor macros.

- `scx_stats` contains endpoint for statistics collection, and can be easily extended. (Here this is unused.)

- New language which is hard for beginners (like me).
  - Traits, lifetimes, ownership is all touched upon by scx package.
  - Async Rust is even harder with its complicated concepts on ownership and lifetimes.

- Not necessarily safe, as it uses unsafe code to interface C pointers.

2. Interlude: Plugins in C

A plugin structure is used to hide away all the complexity of writing anything in eBPF C with restraints.

Essentially, to implement a scheduler policy, we need:

- A `task_struct` to provide the context of the task to be scheduled.

A scheduler policy only need to provide:

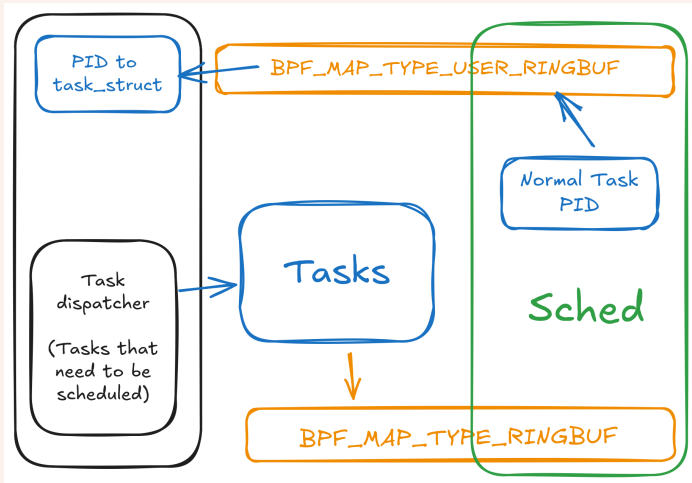- A pid of the process that it decides to run.

Essentially, to implement a scheduler policy, we need:

- A `task_struct` to provide the context of the task to be scheduled.

A scheduler policy only need to provide:

- A pid of the process that it decides to run.

This, combined with the zero-copy `BPF_MAP_TYPE_RINGBUF` and `BPF_MAP_TYPE_USER_RINGBUF` maps, allows us to efficiently pass data between kernel and userspace.

Although we still need to cross boundaries when we need the userspace to schedule processes, but the overhead of crossing the boundary is much smaller.

Still, to schedule a process we still need to cross the boundary twice and call the userspace functions.
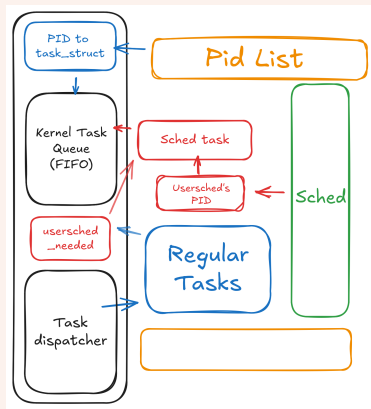
Although we still need to cross boundaries when we need the userspace to schedule processes, but the overhead of crossing the boundary is much smaller.

Still, to schedule a process we still need to cross the boundary twice and call the userspace functions.

To minimize this problem this we can batch the scheduling decisions in the userspace by:

- Using a FIFO queue in eBPF to run scheduling decisions.

- Make the userspace supply an array of $n$ pids to be scheduled.

- Under high load, the ratio of scheduler by actual content is roughly $1/n$.

Another issue is that if the userspace scheduler is not scheduled to run, it cannot make decisions, causing deadlock.
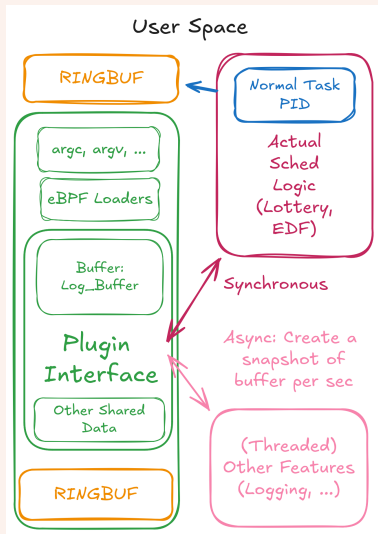


Inspired by scx_userland scheduler, we make these design:

- The userspace will pass its PID to the eBPF scheduler.

- The eBPF scheduler will always push that PID to the scheduling queue rather to userspace.

- Atomic flag is used to notify when there is other task to be scheduled.

In order to observe the behavior of the scheduler, logging is needed.
We may also need other features that has other functions.

Run everything in synchronous, and the scheduler will blow up. I/O
and other things that we need to run is much more expensive then
scheduling.

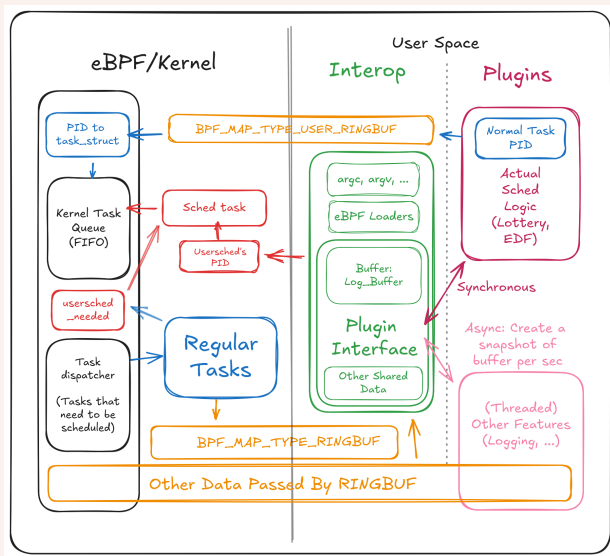As we are already this deep in the rabbit hole we decided to use plugin
and dynamic libraries to allow runtime loading.

- Only one scheduler will be run, and it is the only plugin to run in sync.

- Other plugins have its own thread spanned by the plugin interface.

- The scheduler can log its decisions into a circular buffer.

- After a time interval, the async plugins receive the head and tail to the circular buffer.

While we designed the circular buffer to have 500000 entries to store huge loads, we also introduced a *debug* mode:

- In debug mode, the circular buffer will be alloc'ed 1GB in memory, the exact size of a huge transparent page. This currently holds more than 10 million scheduling decisions.

- The logging plugins will not operate regularly.

- Instead, it will only operate on quit of the scheduler plugin.

- This allows us to have a complete log of the scheduling decisions in extended interval.

Thus we have discussed all the structure relevant to the plugin system, and can turn back to this picture.

```
1   compute_0   Sched: 54130
2   compute_1   Sched: 51487
3   compute_2   Sched: 48143
4   compute_3   Sched: 45032
5   compute_4   Sched: 40560
6   compute_5   Sched: 36630
7   compute_6   Sched: 32557
8   compute_7   Sched: 28394
9   compute_8   Sched: 24328
10  compute_9   Sched: 20886
11  compute_10  Sched: 17503
12  compute_11  Sched: 14884
13  compute_12  Sched: 12187
14  compute_13  Sched: 10212
15  compute_14  Sched: 8246
16  compute_15  Sched: 6665
17  compute_16  Sched: 5374
18  compute_17  Sched: 4367
19  compute_18  Sched: 3435
20  compute_19  Sched: 2815
```

```
1   compute_20 Sched: 2317
2   compute_21 Sched: 1837
3   compute_22 Sched: 1474
4   compute_23 Sched: 1217
5   compute_24 Sched: 931
6   compute_25 Sched: 761
7   compute_26 Sched: 592
8   compute_27 Sched: 453
9   compute_28 Sched: 376
10  compute_29 Sched: 308
11  compute_30 Sched: 244
12  compute_31 Sched: 168
13  compute_32 Sched: 156
14  compute_33 Sched: 129
15  compute_34 Sched: 136
16  compute_35 Sched: 83
17  compute_36 Sched: 71
18  compute_37 Sched: 61
19  compute_38 Sched: 40
20  compute_39 Sched: 39
```

The final structure of the plugin system is good for:

- Simpler implementation of a scheduler policy.

```
1  struct scx_scheduler_plugin_info scx_plugin_info = {
2      .info = {.capability = sched,},
3      .init = fifo_init,
4      .exit = fifo_exit,
5      .enqueue = fifo_enqueue,
6      .dispatch = fifo_dispatch,
7  }
```

- Plugin everything, slightly modified my teammates' pressure tester
  C++ and compiled to a dynamic library, and now the test load is
  also a plugin.

The final structure of the plugin system is flawed for:

- There still lacks a command line capability to specify plugins to load.

- There still lacks a fallback mechanism when no scheduler plugin is loaded. Also bug-prone plugins will stuck the whole computer. *(Pass four null pointers as functions?)*

- The interface(header) exposed to plugins is currently a bit messy and needs refactoring.

That said, these problems exist because of time constraints, and can be solved in future work.

Aside from before problems, there are still some questions:

- In systems with multi CPU it needs fine grained policy and interface.

- How to implement both the core eBPF and the plugin interface for things like:
  - CPU affinity
  - Cache locality
  - NUMA nodes
  - ...

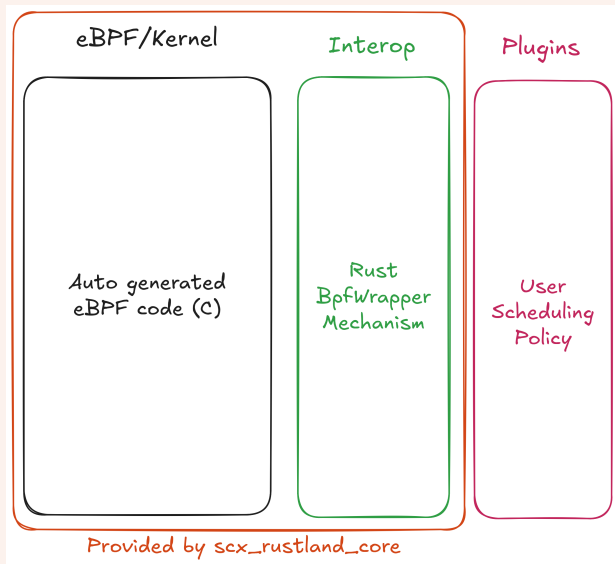Aside from before problems, there are still some questions:

- In systems with multi CPU it needs fine grained policy and interface.

- How to implement both the core eBPF and the plugin interface for things like:
    - CPU affinity
    - Cache locality
    - NUMA nodes
    - ...

- The aim is to have a scalable and flexible scheduler that has a very reasonable default but gives the flexibility for user to implement specific policies.

Maybe it is time to turn back to Rust for help from existing crates?

# 3. Rust Scheduler

We can turn our attention to writing the entire scheduler in Rust.

This uses the crate scx_rustland_core and scx_utils to provide the eBPF bindings and the plugin interface.

The main function of these crates is to auto-generate eBPF code and corresponding Rust bindings to let us focus on writing the scheduling policy only.

eBPF/Kernel

Interop

Plugins

Auto generated
eBPF code (C)

Rust
BpfWrapper
Mechanism

User
Scheduling
Policy

Provided by scx_rustland_core

Instead from the
bottom up like
in the last part,
we start from
the top.

These generated codes can be modified, but we use them more as a reasonable default for our time constraints.

We will introduce it in three aspects, compared to my previous C plugin implementation:

- How I used this as a basis to build the scheduler in Rust.

- How the abstraction is designed in Rust.

- The similarities in implementation of eBPF.

We start with a Scheduler object:

```
1  struct Scheduler<'a> {
2      bpf: BpfScheduler<'a>,   // BPF connector
3      tasks: BTreeSet<Task>,   // tasks ordered by deadline
4      slice_ns: u64,           // Default time slice (in ns)
5  }
```

and a schedule method:

```
1  fn schedule(&mut self) {
2      self.enqueue_all_from_bpf();
3      self.dispatch_task_to_bpf(8);
4      // Notify the dispatcher if there are still
5      // pending tasks to be processed.
6      self.bpf.notify_complete(self.tasks.len() as u64);
7  }
```

# Some Psuedocode

```rust
fn enqueue_all_from_bpf(&mut self) {
    let timestamp = Self::now();
    loop { match self.bpf.dequeue_task() {
        Ok(Some(task_in)) =>
            // Preserve the previous valid deadline
            let deadline = if task_in.vtime != 0 {
                task_in.vtime
            } else { Self::niceness_to_time()};
            self.tasks.insert(Task);

        Ok(None) => ...
        Err(err) => ...
        }
    }
}
```

```
1  // dispatch some tasks to ebpf
2  fn dispatch_task_to_bpf(&mut self, num: u64) {
3      let Some(task) = self.tasks.pop_first() else {
4          break;
5      };
6      let mut dispatched_task =
7          DispatchedTask::new(&task.qtask);
8
9      dispatched_task.slice_ns = ...
10     dispatched_task.vtime = task.deadline;
11     dispatched_task.cpu = ...
12
13     // Send the task to the BPF dispatcher.
14     if self.bpf.dispatch_task(&dispatched_task).is_err()
15         // Dispatching failed: reinsert the task
16         self.tasks.insert(task); break;
17 }
```

```
1   compute_0  Sched: 54034
2   compute_1  Sched: 52112
3   compute_2  Sched: 51831
4   compute_3  Sched: 45332
5   compute_4  Sched: 50409
6   compute_5  Sched: 47692
7   compute_6  Sched: 41756
8   compute_7  Sched: 37856
9   compute_8  Sched: 32610
10  compute_9  Sched: 25936
11  compute_10 Sched: 21448
12  compute_11 Sched: 16813
13  compute_12 Sched: 13327
14  compute_13 Sched: 10580
15  compute_14 Sched: 8318
16  compute_15 Sched: 6160
17  compute_16 Sched: 4944
18  compute_17 Sched: 4523
19  compute_18 Sched: 3201
20  compute_19 Sched: 3052
```

```
1   compute_20 Sched: 2473
2   compute_21 Sched: 2018
3   compute_22 Sched: 1428
4   compute_23 Sched: 1405
5   compute_24 Sched: 1129
6   compute_25 Sched: 933
7   compute_26 Sched: 766
8   compute_27 Sched: 642
9   compute_28 Sched: 527
10  compute_29 Sched: 402
11  compute_30 Sched: 396
12  compute_31 Sched: 273
13  compute_32 Sched: 281
14  compute_33 Sched: 186
15  compute_34 Sched: 154
16  compute_35 Sched: 171
17  compute_36 Sched: 123
18  compute_37 Sched: 105
19  compute_38 Sched: 105
20  compute_39 Sched: 69
```

Now we can turn to the implementation of scx_rustland_core's bpf.rs and see how it abstracts away the eBPF code.
Similarly, it also uses objects to represent the BPF scheduler.

```rust
pub struct BpfScheduler<'cb> {
    pub skel: BpfSkel<'cb>,
        // low level BPF skeleton
    shutdown: Arc<AtomicBool>,
        // Atomic flag for shutdown
    queued: libbpf_rs::RingBuffer<'cb>,
    dispatched: libbpf_rs::UserRingBuffer,
    struct_ops: Option<libbpf_rs::Link>,
        // Low level BPF methods
}
```

**Task Dequeue Abstraction**

```
1 pub fn dequeue_task(&mut self) -> Result<Option<QueuedTask>, i32>
```

**Task Dispatch Abstraction**

```
1 pub fn dispatch_task(&mut self, task: &DispatchedTask)
2     -> Result<(), libbpf_rs::Error>
```

**Type-Safe Counter Access**

```
1  pub fn nr_queued_mut(&mut self) -> &mut u64
2  pub fn nr_scheduled_mut(&mut self) -> &mut u64
3  pub fn nr_user_dispatches_mut(&mut self) -> &mut u64
4  // ... other counter methods
```

- Type-safe counter access without manual offset calculations

- Direct mapping to BPF memory regions

- Eliminates pointer arithmetic and memory safety risks

- Automatic synchronization with kernel space

**Completion Notification Abstraction**

```
1  pub fn notify_complete(&mut self, nr_pending: u64)
```

- Notifies BPF that userspace scheduler completed scheduling cycle

- Updates count of pending tasks

- Prevents unnecessary BPF busy-looping

Now we can turn our attention to the eBPF implementation details, which is also in C.

The eBPF implementation in here, while much more detailed, on a higher perspective is similar to that in C plugin:

- Both uses two ringbuf to pass data between kernel and userspace.

- Both mostly passes structs to userspace and gets information back.

Multi-Level Prevention Strategy

1. Atomic Flag Notification:

```
volatile u32 usersched_needed;        // atomic flag
set_usersched_needed();               // atomic OR
test_and_clear_usersched_needed();    // atomic AND
```

2. Heartbeat Timer Mechanism:

```
usersched_timer_fn();  // checks every USERSCHED_TIMER_NS (1 sec)
if (scheduler inactive > 1 sec) {
    set_usersched_needed();  // wake scheduler
    kick CPU;
}
```

# Differences: Wakeup Logic

3. Dedicated DSQ for Scheduler:

```
1  #define SCHED_DSQ (MAX_CPUS + 1)  // dedicated queue
2  // In rustland_enqueue(): place scheduler task in SCHED_DSQ
3  // In rustland_dispatch(): only dispatch from SCHED_DSQ if needed
```

4. Comprehensive Wakeup Check:

```
1  bool usersched_has_pending_tasks() {
2      return atomic_flag(usersched_needed) ||
3             nr_scheduled > 0 ||
4             bpf_ringbuf_query(&queued, BPF_RB_AVAIL_DATA);
5  }
```

Result: Scheduler only wakes when there is actual work; reduces unnecessary context switches.

- It is still very dependent on C eBPF to get the **mechanism** running in kernel space.

- But through this plugin method we can make actual **policy** in any way (language, method, framework) we like.

- The ecosystem/build system of Rust makes it possible to import existing frameworks quickly and easily, which is the core advantages compared to C.

- And it has second best ecosystem and community only to C, an advantage over other languages.

- Along with other modern features it is easy to play around, build and extend without reinventing the wheel for everything like in C.
  - Quickly try new policies
  - Use different data structures
  - Different logging methods
  - . . .

- That said Rust is indeed hard to learn.

- To understand eBPF it is better to at least write a working plugin architecture from the ground up like what I did in C plugin, and then turn to higher-level abstractions.

**Thank you!**

We tried our best to follow agile development practices, communication guidelines, etc.
We also used git hooks for commit checking to ensure atomic commit.
We only broke a few times mainly due to copying infrastructure code around.

The contribution on gitea is not proportional, as I did a lot of work that was not strictly required.

Some code was copied over and over in different parts of the project, especially in the C plugin part.

Rust configuration and copied build scripts also took a some commits.