

Readme

本文档为本次“操作系统大赛”的概括性文档，包含了本次大赛参赛作品的一些基本信息和TO-DO List。具体的教程文档，存放在对应的文件夹中。

因为指令集移植工作来不及做，所以部分代码仍然是x86风格

一、实验模块设计

实验节点要求满足模块化的特点。我们整个实验分为以下几个模块：

1、最小化内核：

本部分要求对操作系统的启动过程有足够的了解，并且能够编写并启动一个能够打印Hello World的最小化内核。包括了如下的步骤：

- 搭建开发环境。可以在Ubuntu等Linux发行版中进行。也可以在Windows操作系统中，配合WSLG等工具进行开发。
- 安装相应的开发工具，例如VSCode，rustc等。
- 了解rust的SBI，并学习借助sbi载入操作系统的方法。K210开发板中存在uboot。（待填坑）
- 利用如上的工具和机制，通过加载ELF文件的方式，进入操作系统内核。并且在内核中打印Hello World。这样可以大大减少汇编语言的使用。
- 如下是进入OS的主函数，这是移植之前的，所以使用的是Rust中与uefi有关的包。移植时将使用sbi完成类似的操作。这样做的一大好处就是可以最大程度减少汇编语言的使用。

解释待填坑，这一块我实在不太懂了（by brh）

```
1  #[entry]
2  fn efi_main(image: uefi::Handle, mut system_table: SystemTable<Boot>) ->
   Status {
3      uefi_services::init(&mut system_table).expect("Failed to initialize
   utilities");
4
5      info!("Running UEFI bootloader...");
6
7      let bs = system_table.boot_services();
8      let config = {
9          let mut file = open_file(bs, CONFIG_PATH);
10         let buf = load_file(bs, &mut file);
11         config::Config::parse(buf)
12     };
13
14     let graphic_info = init_graphic(bs);
15     // info!("config: {:#x?}", config);
16
17     let acpi2_addr = system_table
18         .config_table()
19         .iter()
20         .find(|entry| entry.guid == ACPI2_GUID)
21         .expect("failed to find ACPI 2 RSDP")
22         .address;
```

```

23     info!("ACPI2: {:?}", acpi2_addr);
24
25     let elf = {
26         let mut file = open_file(bs, config.kernel_path);
27         let buf = load_file(bs, &mut file);
28         ElfFile::new(buf).expect("failed to parse ELF")
29     };
30     unsafe {
31         ENTRY = elf.header.pt2.entry_point() as usize;
32     }
33
34     let max_mmap_size = system_table.boot_services().memory_map_size();
35     let mmap_storage = Box::leak(
36         vec![0; max_mmap_size.map_size + 10 *
max_mmap_size.entry_size].into_boxed_slice()
37     );
38     let mmap_iter = system_table
39         .boot_services()
40         .memory_map(mmap_storage)
41         .expect("Failed to get memory map")
42         .1;
43     let max_phys_addr = mmap_iter
44         .map(|m| m.phys_start + m.page_count * 0x1000)
45         .max()
46         .unwrap()
47         .max(0x1_0000_0000); // include IOAPIC MMIO area
48
49     let mut page_table = current_page_table();
50     // root page table is readonly
51     // disable write protect
52     unsafe {
53         Cr0::update(|f| f.remove(Cr0Flags::WRITE_PROTECT));
54         Efer::update(|f| f.insert(EferFlags::NO_EXECUTE_ENABLE));
55     }
56
57     elf::map_elf(&elf, &mut page_table, &mut UEFIFrameAllocator(bs))
58         .expect("Failed to map ELF");
59
60     elf::map_range(
61         config.kernel_stack_address,
62         config.kernel_stack_size,
63         &mut page_table,
64         &mut UEFIFrameAllocator(bs),
65         false
66     ).expect("Failed to map stack");
67
68     elf::map_physical_memory(
69         config.physical_memory_offset,
70         max_phys_addr,
71         &mut page_table,
72         &mut UEFIFrameAllocator(bs),
73     );
74
75     // recover write protect
76     unsafe {

```

```

77         Cr0::update(|f| f.insert(Cr0Flags::WRITE_PROTECT));
78     }
79
80     // FIXME: multi-core
81     // All application processors will be shutdown after ExitBootService.
82     // Disable now.
83     // start_aps(bs);
84
85     // for i in 0..5 {
86     //     info!("waiting for next stage... {}", 5 - i);
87     //     bs.stall(100_000);
88     // }
89
90     info!("Exiting boot services...");
91
92     let (rt, mmap_iter) = system_table
93         .exit_boot_services(image, mmap_storage)
94         .expect("Failed to exit boot services");
95     // NOTE: alloc & log can no longer be used
96
97     // construct BootInfo
98     let bootinfo = BootInfo {
99         memory_map: mmap_iter.copied().collect(),
100         physical_memory_offset: config.physical_memory_offset,
101         graphic_info,
102         system_table: rt,
103     };
104     let stacktop = config.kernel_stack_address + config.kernel_stack_size *
105     0x1000;
106     unsafe {
107         jump_to_entry(&bootinfo, stacktop);
108     }
109 }

```

2、控制台、日志和调试：

本部分是了解一些常见的调试方式，使得后续的实验轻松简单。

- 利用串口输出作为调试。这一步需要手动实现println等功能。在使用log crate的前提下，
- 利用VGA输出进行调试（可选）
- 编译时附带调试信息。编译时除了编译为Debug版本和Release版本，还可以选择编译为Release with debug info的版本，也即可以对Release版本进行调试。
- vscode调试。本部分皓宇记得填坑，着重介绍
- gdb调试：GDB调试在操作系统的编写中是相当重要的。借助于插件pwndbg，可以使得debug变得容易很多。
- LLDB：内置于XCode中的调试器，可以使得使用Mac的同学轻松地Debug。其他的平台也可以安装使用。

重点展示这一部分，皓宇记得写一下，这是亮点

3、内存管理（不含缺页中断）

- RISC-V SV39内存管理：对于64位的RISC-V架构，多使用三级页表，支持39位虚拟地址，也即每个地址池理论上最多支持512GB的内存（事实上被切分为两个256GB）。物理地址显然为64位。虚拟地址通过MMU转换为物理地址。分页模式选择和页表基地址保存在SARF寄存器中。
- 页面分配：使用BitMap管理资源的分配。先分配连续的虚拟页，再给每个虚拟页分配对应的物理页。
- 动态内存分配：在虚存的 `0xFFFFFFFF8000000000` 开始分配一个 32MB 的堆，从 Bootloader 传来的 mmap 中可以拿到可用的内存区域，在帧分配器中可以将它们切成 4KiB 的块，在我们需要的时候进行分配。同时，我们也可以给帧分配器附带一个动态数组，用于存储被释放的物理帧，以供再次使用

4、中断与陷阱

- RISCV的中断依赖于mtvec和mcause寄存器
- mtvec寄存器存放了中断处理程序的基地址和中断处理函数的寻址模式。我们选择直接寻址，并用模式匹配的方法去匹配到合适的函数。
- mcause则记录了中断发生的原因。
- Rust中的cpu crate有TrapFrame结构体，记录了一次中断的整型以及浮点寄存器、MMU、中断堆栈地址和内核线程编号。
- 对于不同类型的中断，我们需要单独编写不同的中断处理程序进行处理。
- 对于外部中断，则有赖于PLIC进行处理。这有赖于mie寄存器中的meie位

5、内核线程：

- 前面已经实现了动态内存分配，我们就可以为每个进程分配一块内存空间作为PCB。PCB中当然至少应该包括栈信息、帧信息、pid、计数、页目录基地址和状态。
- 要新生成一个线程，需要申请栈帧和页目录表，然后将新创建的内核线程加入到队列中

```
1 pub fn spawn_kernel_thread(entry: fn() -> !, name: String, data:
   Option<ProcessData>) {
2     x86_64::instructions::interrupts::without_interrupts(|| {
3         let entry = VirtAddr::new(entry as u64);
4
5         let stack = get_frame_alloc_for_sure().allocate_frame()
6             .expect("Failed to allocate stack for kernel thread");
7
8         let stack_top = VirtAddr::new(physical_to_virtual(
9             stack.start_address().as_u64()) + FRAME_SIZE);
10
11         let mut manager = get_process_manager_for_sure();
12         manager.spawn_kernel_thread(entry, stack_top, name, ProcessId(0),
13             data);
14     });
15 }
16
17 impl ProcessManager {
18     pub fn spawn_kernel_thread(
19         &mut self,
20         entry: VirtAddr,
21         stack_top: VirtAddr,
22         name: String,
23         parent: ProcessId,
24         proc_data: Option<ProcessData>,
25     ) -> ProcessId {
```

```

25     let mut p = Process::new(
26         &mut *crate::memory::get_frame_alloc_for_sure(),
27         name,
28         parent,
29         self.get_kernel_page_table(),
30         proc_data,
31     );
32     p.pause();
33     p.init_stack_frame(entry, stack_top);
34     info!("Spawn process: {}#{}", p.name(), p.pid());
35     let pid = p.pid();
36     self.processes.push(p);
37     pid
38 }
39 }

```

- 调度方法是时间片轮转调度。当记录到一定次数的时钟中断时，或者进程结束时，就启动调度程序。方法在于先保存目前栈帧的信息，然后修改PC寄存器和栈指针，恢复栈帧信息，就切换到了另一个进程执行。

- 当然其他的进程调度算法也是可用的

```

1  pub fn switch(regs: &mut Registers, sf: &mut InterruptStackFrame) {
2      x86_64::instructions::interrupts::without_interrupts(|| {
3          let mut manager = get_process_manager_for_sure();
4
5          manager.save_current(regs, sf);
6          manager.switch_next(regs, sf);
7      });
8  }
9
10 impl ProcessManager {
11     pub fn save_current(&mut self, regs: &mut Registers, sf: &mut
12     InterruptStackFrame) {
13         let current = self.current_mut();
14         if current.is_running() {
15             current.tick();
16             current.save(regs, sf);
17         }
18         // trace!("Paused process #{}", self.cur_pid);
19     }
20
21     pub fn switch_next(&mut self, regs: &mut Registers, sf: &mut
22     InterruptStackFrame) {
23         let pos = self.get_next_pos();
24         let p = &mut self.processes[pos];
25
26         // trace!("Next process {} #{}", p.name(), p.pid());
27         if p.pid() == self.cur_pid {
28             // the next process to be resumed is the same as the current one
29             p.resume();
30         } else {
31             // switch to next process
32             p.restore(regs, sf);
33             self.cur_pid = p.pid();
34         }
35     }
36 }

```

```

32     }
33 }
34 }

```

- 进程间数据访问的互斥和同步可以通过互斥锁、信号量和管程实现，这些都不复杂。
- 进程间还可以通过管道通信。

6、I/O的处理

- 这里主要指的是获取键盘和串口的输入
- 输入的实现一般会有一个循环队列作为缓存，用来保存输入的键值，在需要输出或使用的时候取出，这样就可以保证输入的正常顺序。
- 对于键盘和串口的输入，可以提供一套统一的接口，也即键盘和串口输入的数据统统放进同一个队列，而需要输入的地方只需要从队列中取用、并在队列为空时等待即可。
- `InputStream` 结构体的作用是，初始化一个输入队列、并承担弹出字符的作用。

```

1  pub struct InputStream;
2
3  impl InputStream {
4      pub fn new() -> Self {
5          init_INPUT_QUEUE(ArrayQueue::new(DEFAULT_BUF_SIZE));
6          info!("Input stream Initialized.");
7          Self
8      }
9  }
10
11 impl Stream for InputStream {
12     type Item = DecodedKey;
13
14     fn poll_next(self: Pin<&mut Self>, cx: &mut Context) ->
15     Poll<Option<Self::Item>> {
16         let queue = get_input_queue_for_sure();
17         if let Some(key) = queue.pop() {
18             Poll::Ready(Some(key))
19         } else {
20             INPUT_WAKER.register(&cx.waker());
21             match queue.pop() {
22                 Some(key) => {
23                     INPUT_WAKER.take();
24                     Poll::Ready(Some(key))
25                 }
26                 None => Poll::Pending,
27             }
28         }
29     }
30 }

```

- 需要注意的是，更多时候我们的输入并不是简单的字符，可能一串字节流在解析后是一个控制字符等等，因此我们也需要一个全局的键盘解释器，这样一个静态解释器也需要从两方同时获取输入。如果可以解析为Unicode字符那就直接输出，否则按照调试模式输出。`DecodedKey` 这个枚举类型可以直接调用现有的库。

```

1 pub async fn get_key() {
2     let mut input = InputStream::new();
3     while let Some(key) = input.next().await {
4         match key {
5             DecodedKey::Unicode(c) => print!("{}", c),
6             DecodedKey::RawKey(k) => print!("{:?}", k),
7         }
8     }
9 }

```

7、文件系统与页面交换

- 磁盘是一个块设备，可能含有多个分区，我们采用MBR分区表对其进行分区。分区表内部包含了磁盘大小、各分区大小等信息。MBR分区表支持至多四个物理分区。
- 每个分区可以有不同的文件系统。我们先实现了一个FAT16的文件系统，这个文件系统的实现相对简单，可以让做实验的同学不太困难地写出来。如果同学学有余力可以尝试更复杂的文件系统
-

8、系统调用

- 在RISC-V中，调用系统调用需要使用ecall指令，可以调用8号内部中断
- 我们可以在寄存器a0中存放系统调用号，a1~a7放参数。这样去使用系统调用。使用方法大致类似于C语言。
- 考虑到Rust强大的模式匹配能力，我们不需要再去编写系统调用表。
- 我们可以按照一定的规范（如POSIX），编写系统调用。
- 以下的代码仍然是x86-64下触发0x80中断的。但是两种架构下调用系统调用的思路大同小异。同样可以使用模式匹配来找到系统调用函数。

```

1 pub fn dispatcher(regs: &mut Registers, sf: &mut InterruptStackFrame) {
2     let args = super::syscall::SyscallArgs::new(
3         Syscall::try_from(regs.rax as u8).unwrap(),
4         regs.rdi,
5         regs.rsi,
6         regs.rdx
7     );
8
9     trace!("{}", args);
10
11     match args.syscall {
12         // path: &str (arg0 as *const u8, arg1 as len) -> pid: u16
13         Syscall::Spawn => regs.set_rax(spawn_process(&args)),
14         // pid: arg0 as u16
15         Syscall::Exit  => exit_process(&args, regs, sf),
16         // fd: arg0 as u8, buf: &[u8] (arg1 as *const u8, arg2 as len)
17         Syscall::Read  => regs.set_rax(sys_read(&args)),
18         // fd: arg0 as u8, buf: &[u8] (arg1 as *const u8, arg2 as len)
19         Syscall::write => regs.set_rax(sys_write(&args)),
20         // path: &str (arg0 as *const u8, arg1 as len), mode: arg2 as u8 ->
        fd: u8
21         Syscall::Open  => regs.set_rax(sys_open(&args)),
22         // fd: arg0 as u8 -> success: bool

```

```

23     Syscall::close          => regs.set_rax(sys_close(&args)),
24     // None
25     Syscall::Stat           => list_process(),
26     // None -> time: usize
27     Syscall::Time           => regs.set_rax(sys_clock() as usize),
28     // path: &str (arg0 as *const u8, arg1 as len)
29     Syscall::ListDir        => list_dir(&args),
30     // layout: arg0 as *const Layout -> ptr: *mut u8
31     Syscall::Allocate       => regs.set_rax(sys_allocate(&args)),
32     // ptr: arg0 as *mut u8
33     Syscall::Deallocate     => sys_deallocate(&args),
34     // x: arg0 as i32, y: arg1 as i32, color: arg2 as u32
35     Syscall::Draw           => sys_draw(&args),
36     // pid: arg0 as u16 -> status: isize
37     Syscall::WaitPid        => regs.set_rax(sys_wait_pid(&args)),
38     // None -> pid: u16
39     Syscall::GetPid         => regs.set_rax(sys_get_pid() as usize),
40     // None -> pid: u16 (diff from parent and child)
41     Syscall::Fork           => sys_fork(regs, sf),
42     // pid: arg0 as u16
43     Syscall::Kill           => sys_kill(&args, regs, sf),
44     // op: u8, key: u32, val: usize -> ret: any
45     Syscall::Sem            => sys_sem(&args, regs, sf),
46     // None
47     Syscall::None           => {}
48 }
49 }

```

9、程序加载与用户进程

- 此时考虑从磁盘中加载程序到内存，这里一般指的是加载ELF文件
- ELF文件一般包括如下的部分：ELF Header、Program Headers和可执行程序部分
- 使用相应的结构体解析ELF Header、Program Headers。如果发现执行权限和架构都合适，那么就可以执行。
- 随后可以构建PCB，执行用户进程。
- 用户进程中的动态内存分配，可以在内核中创建一个共用的动态内存分配器
- ELF文件亦有相关的crate可供使用。


```

1 pub fn load_elf(
2     elf: &ElfFile,
3     physical_offset: u64,
4     page_table: &mut impl Mapper<Size4KiB>,
5     frame_allocator: &mut impl FrameAllocator<Size4KiB>,
6     user_access: bool,
7 ) -> Result<Vec<PageRangeInclusive>, MapToError<Size4KiB>> {
8     trace!("Loading ELF file...{:?}", elf.input.as_ptr());
9     elf.program_iter()
10        .filter(|segment| segment.get_type().unwrap() == program::Type::Load)
11        .map(|segment| load_segment(elf, physical_offset, &segment,
12        page_table, frame_allocator, user_access))
13        .collect()
14 }

```

10、shell

- shell的本质是，使用fork和exec运行一些特定的程序。Linux的shell支持管道，可以将上一次执行的结果作为下一次执行的输入
- 对于简单的shell实现，可以直接解析字符串，然后运行对应的函数来实现功能。

```

1 fn main() -> usize {
2     let mut root_dir = String::from("/APP/");
3     println!("          <<< welcome to GGOS shell >>>          ");
4     println!("          type `help` for help");
5     loop {
6         print!("[{}] ", root_dir);
7         let input = stdin().read_line();
8         let line: Vec<&str> = input.trim().split(' ').collect();
9         match line[0] {
10             "exit" => break,
11             "ps" => sys_stat(),
12             "ls" => sys_list_dir(root_dir.as_str()),
13             "cat" => {
14                 if line.len() < 2 {
15                     println!("Usage: cat <file>");
16                     continue;
17                 }
18                 services::cat(line[1], root_dir.as_str());
19             }
20             "cd" => {
21                 if line.len() < 2 {
22                     println!("Usage: cd <dir>");
23                     continue;
24                 }
25                 services::cd(line[1], &mut root_dir);
26             }
27             "exec" => {
28                 if line.len() < 2 {
29                     println!("Usage: exec <file>");
30                 }
31             }
32         }
33     }
34 }

```

```

32         continue;
33     }
34
35     services::exec(line[1], root_dir.as_str());
36 }
37 "nohup" => {
38     if line.len() < 2 {
39         println!("Usage: nohup <file>");
40         continue;
41     }
42
43     services::nohup(line[1], root_dir.as_str());
44 }
45 "kill" => {
46     if line.len() < 2 {
47         println!("Usage: kill <pid>");
48         continue;
49     }
50     let pid = line[1].to_string().parse::<u16>();
51
52     if pid.is_err() {
53         errln!("Cannot parse pid");
54         continue;
55     }
56
57     services::kill(pid.unwrap());
58 }
59 "help" => print!("{}", consts::help_text()),
60 _ => println!("[=] you said \"{}\"", input),
61 }
62 }
63
64 0
65 }

```

11、多核支持

待实现

12、图形界面

待实现

13、网络与套接字编程

待实现

14、指令集移植

- 因为我们的这个项目是基于一个自己编写的X86-64操作系统，我们在做这一步的时候其实是把x86-64移植到RISC-V，然后教程反过来写就可以了。
- 要进行指令集移植，我们必须在如下方面进行更改：
 - 寄存器结构不同：riscv提供了大量的通用寄存器，但是x86-64的通用寄存器较少（虽然比x86还是多很多的），专用寄存器很多。而且一些特殊寄存器二者完全不同。
 - 汇编指令不同：这是显然的。主要注意一点：x86-64的外设是独立编址的，但是riscv的外设是统一编址的。这将涉及到所有的外设，比如硬盘和键盘。在这里，x86-64架构在rust中有很多可用的crate帮助减少汇编的使用，但是对于riscv架构，支持相对缺乏，可能使用汇编的地方相对较多。
 - 启动：uefi和rust-sbi的启动有所不同。**待填坑**
 - 页表：RISC-V多使用SV39三级页表（也有SV48四级页表可供使用），x86-64多使用四级页表。
 - 中断和系统调用：二者调用中断和系统调用的方式也有很大区别。对于x86-64需要借助中断描述符表但是riscv是不需要的

各个模块的依赖关系如下：

(一张图片)

理由（若干字）

二、已经实现的部分：

- 本组已经设计了一个基于x86-64和uefi的操作系统。目前九点绝大多数已经实现。文件系统的写入功能和交换文件还有待填坑。
-