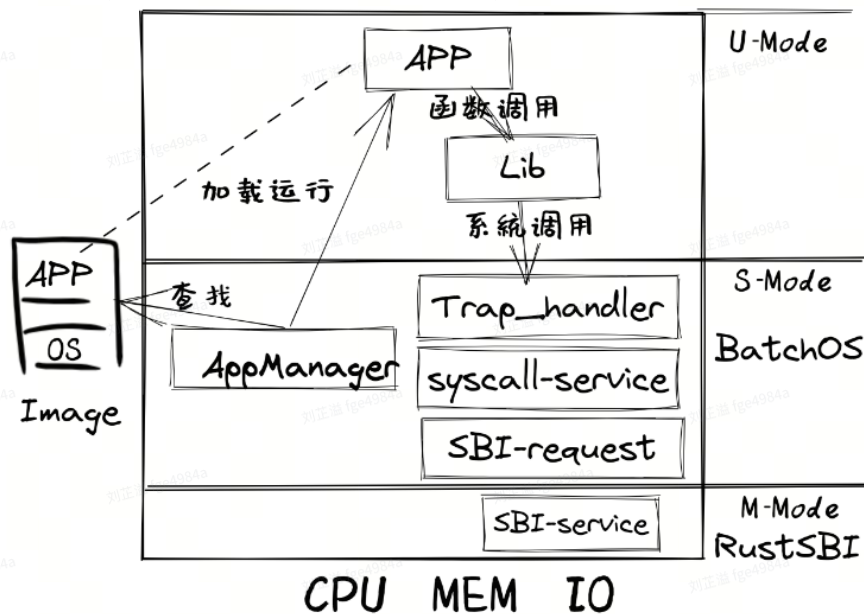


BatchOS

BatchOS



- 实现批处理程序功能的OS
- APP与OS隔离
- 自动加载并运行多个程序

特权级机制

- `ecall` 具有用户态到内核态的执行环境切换能力的函数调用指令；
- `sret`：具有内核态到用户态的执行环境切换能力的函数返回指令。
- 首先，操作系统需要提供相应的功能代码，能在执行 `sret` 前准备和恢复用户态执行应用程序的上下文。其次，在应用程序调用 `ecall` 指令后，能够检查应用程序的系统调用参数，确保参数不会破坏操作系统。

RISC-V异常

Interrupt	Exception Code	Description
0	0	Instruction address misaligned
0	1	Instruction access fault
0	2	Illegal instruction
0	3	Breakpoint
0	4	Load address misaligned
0	5	Load access fault
0	6	Store/AMO address misaligned
0	7	Store/AMO access fault
0	8	Environment call from U-mode
0	9	Environment call from S-mode
0	11	Environment call from M-mode
0	12	Instruction page fault
0	13	Load page fault
0	15	Store/AMO page fault

RISC-V S模式特权指令

- sret: 从S模式返回U模式
- wfi: 处理器在空闲时进入低功耗状态等待终端
- sfence.vma: 刷新TLB缓存
- 访问S模式CSR指令: 改变系统状态

控制状态寄存器

- sstatus: SPP字段给出Trap发生前CPU的特权级
- sepc: 记录异常发生前执行的最后一条指令的地址

- `scause`: 描述Trap的原因
- `stval`: 给出Trap的附加信息
- `stvec`: 控制Trap处理代码的入口地址

硬件切换的硬件控制机制

ecall

- `sstatus` 中的 `SPP` 字段切换到CPU当前特权级
- `sepc` 修改为Trap处理完成后默认执行的下一条指令的地址
- `scause\stval` 修改成Trap原因和Trap额外信息
- CPU跳转到 `stvec` 设置的Trap处理入口函数，设置特权级为S

sret

- CPU按照 `sstatus` 设置特权级
- CPU跳转到 `sepc` 指向的地址，然后继续执行

用户库

- 使用 Rust 的宏将其函数符号 `main` 标志为弱链接。这样在最后链接的时候，虽然在 `lib.rs` 和 `bin` 目录下的某个应用程序都有 `main` 符号，但由于 `lib.rs` 中的 `main` 符号是弱链接，链接器会使用 `bin` 目录下的应用主逻辑作为 `main`
- 这里我们主要是进行某种程度上的保护，如果在 `bin` 目录下找不到任何 `main`，那么编译也能够通过，但会在运行时报错。

```
1  #[linkage="weak"]
2  #[no_mangle]
3  fn main()->i32 {
4      panic!("Cannot find main!");
5  }
```

- `#![feature(linkage)]` 支持链接操作

系统调用

- `&[u8]` 切片类型来描述缓冲区，这是一个**胖指针** (Fat Pointer)，里面既包含缓冲区的起始地址，还包含缓冲区的长度。

```
1  // user/src/syscall.rs
2  use core::arch::asm;
```

```

3 fn syscall(id: usize, args: [usize; 3]) -> isize {
4     let mut ret: isize;
5     unsafe {
6         asm!(
7             "ecall",
8             inlateout("x10") args[0] => ret,
9             in("x11") args[1],
10            in("x12") args[2],
11            in("x17") id
12        );
13    }
14    ret
15 }
16 const SYSCALL_WRITE: usize = 64;
17 const SYSCALL_EXIT: usize = 93;
18
19 pub fn sys_write(fd: usize, buffer: &[u8]) -> isize {
20     syscall(SYSCALL_WRITE, [fd, buffer.as_ptr() as usize, buffer.len()])
21 }
22
23 pub fn sys_exit(xstate: i32) -> isize {
24     syscall(SYSCALL_EXIT, [xstate as usize, 0, 0])
25 }

```

加载不同的APP

AppManager

- 在 `RefCell` 的基础上再封装一个 `UPSafeCell`，它名字的含义是：允许我们在单核上安全使用可变全局变量。
- 当我们要访问数据时（无论读还是写），需要首先调用 `exclusive_access` 获得数据的可变借用标记，通过它可以完成数据的读写，在操作完成之后我们需要销毁这个标记，此后才能开始对该数据的下一次访问

```

1 // os/src/sync/up.rs
2
3 pub struct UPSafeCell<T> {
4     /// inner data
5     inner: RefCell<T>,
6 }
7
8 // unsafe向编译器保证只在单核上进行操作
9 unsafe impl<T> Sync for UPSafeCell<T> {}
10

```

```

11 impl<T> UPSafeCell<T> {
12     /// User is responsible to guarantee that inner struct is only used in
13     /// uniprocessor.
14     pub unsafe fn new(value: T) -> Self {
15         Self { inner: RefCell::new(value) }
16     }
17     /// Panic if the data has been borrowed.
18     pub fn exclusive_access(&self) -> RefMut<'_, T> {
19         self.inner.borrow_mut()
20     }
21 }

```

New

- 使用 `core::slice::from_raw_parts` 将指针解释为 `&[usize]` 切片;
- 使用 `copy_from_slice` 将切片上的元素复制到 `app_start` 上

```

1 static ref APP_MANAGER: UPSafeCell<AppManager> = unsafe {
2     UPSafeCell::new({
3         extern "C" {
4             fn _num_app();
5         }
6         let num_app_ptr = _num_app as usize as *const usize;
7         let num_app = num_app_ptr.read_volatile();
8         let mut app_start: [usize; MAX_APP_NUM + 1] = [0; MAX_APP_NUM + 1];
9         let app_start_raw: &[usize] =
10             core::slice::from_raw_parts(num_app_ptr.add(1), num_app + 1);
11         app_start[..=num_app].copy_from_slice(app_start_raw);
12         AppManager {
13             num_app,
14             current_app: 0,
15             app_start,
16         }
17     })
18 };

```

load_app

- CPU用存在指令缓存，使用load_app加载新的程序，需要让OS知道取指内存的变化
- OS 在这里必须使用取指屏障指令 `fence.i`，它的功能是保证在它之后的取指过程必须能够看到在它之前的所有对于取指内存区域的修改

```

1 unsafe fn load_app(&self, app_id: usize) {

```

```

2     if app_id >= self.num_app {
3         println!("All applications completed!");
4         shutdown(false);
5     }
6     println!("[kernel] Loading app_{}", app_id);
7     // clear app area
8     core::slice::from_raw_parts_mut(APP_BASE_ADDRESS as *mut u8,
APP_SIZE_LIMIT).fill(0);
9     // 将指针看作切片，使用拷贝实现App切换
10    let app_src = core::slice::from_raw_parts(
11        self.app_start[app_id] as *const u8,
12        self.app_start[app_id + 1] - self.app_start[app_id],
13    );
14    let app_dst = core::slice::from_raw_parts_mut(APP_BASE_ADDRESS as *mut
u8, app_src.len());
15    app_dst.copy_from_slice(app_src);
16    // Memory fence about fetching the instruction memory
17    // It is guaranteed that a subsequent instruction fetch must
18    // observes all previous writes to the instruction memory.
19    // Therefore, fence.i must be executed after we have loaded
20    // the code of the next app into the instruction memory.
21    // See also: riscv non-priv spec chapter 3, 'Zifencei' extension.
22    asm!("fence.i");
23 }

```

run_next_app

- 先将一个上下文压入内核栈中；在 `__restore` 中更新 `sscrath` 指针指向内核栈栈顶
- 如果发生系统调用，从 `__alltraps` 开始执行

```

1  /// run next app
2  pub fn run_next_app() -> ! {
3      let mut app_manager = APP_MANAGER.exclusive_access();
4      let current_app = app_manager.get_current_app();
5      unsafe {
6          app_manager.load_app(current_app);
7      }
8      app_manager.move_to_next_app();
9      drop(app_manager);
10     // before this we have to drop local variables related to resources
manually
11     // and release the resources
12     extern "C" {
13         fn __restore(cx_addr: usize);
14     }

```

```

15     unsafe {
16         __restore(KERNEL_STACK.push_context(TrapContext::app_init_context(
17             APP_BASE_ADDRESS,
18             USER_STACK.get_sp(),
19             )) as *const _ as usize);
20     }
21     panic!("Unreachable in batch::run_current_app!");
22 }

```

Trap管理

TrapContext

- Trap 发生时需要保存的物理资源内容，包括32个通用寄存器、sstatus以及sepc
- 对于 CSR 而言，我们知道进入 Trap 的时候，硬件会立即覆盖掉 `scause/stval/sstatus/sepc` 的全部或是其中一部分。
 - `scause/stval` 的情况是：它总是在 Trap 处理的第一时间就被使用或者是在其他地方保存下来了，因此它没有被修改并造成不良影响的风险。
 - 而对于 `sstatus/sepc` 而言，它们会在 Trap 处理的全程有意义（在 Trap 控制流最后 `sret` 的时候还用到了它们），而且确实会出现 Trap 嵌套的情况使得它们的值被覆盖掉。所以我们需要将它们也一起保存下来，并在 `sret` 之前恢复原样。

```

1 // os/src/trap/context.rs
2
3 #[repr(C)]
4 pub struct TrapContext {
5     pub x: [usize; 32],
6     pub sstatus: Sstatus,
7     pub sepc: usize,
8 }

```

TrapContext的保存与恢复

- 首先通过 `__alltraps` 将 Trap 上下文保存在内核栈上，然后跳转到使用 Rust 编写的 `trap_handler` 函数完成 Trap 分发及处理。当 `trap_handler` 返回之后，使用 `__restore` 从保存在内核栈上的 Trap 上下文恢复寄存器。最后通过一条 `sret` 指令回到应用程序执行。

```

1 // os/src/trap/mod.rs
2
3 global_asm!(include_str!("trap.S"));

```

```

4
5 pub fn init() {
6     extern "C" { fn __alltraps(); }
7     unsafe {
8         stvec::write(__alltraps as usize, TrapMode::Direct);
9     }
10 }

```

__alltraps

- `sscratch` 在 `__restore` 中设置为内核栈栈顶, `run_next_app` 从 `__restore` 先执行

```

1 # os/src/trap/trap.S
2
3 .macro SAVE_GP n
4     sd x\n, \n*8(sp)
5 .endm
6
7 .align 2
8 __alltraps:
9     csrrw sp, sscratch, sp # exchange sp and sscratch(point to kernel stack)
10    # now sp->kernel stack, sscratch->user stack
11    # allocate a TrapContext on kernel stack
12    addi sp, sp, -34*8
13    # save general-purpose registers
14    sd x1, 1*8(sp)
15    # skip sp(x2), we will save it later
16    sd x3, 3*8(sp)
17    # skip tp(x4), application does not use it
18    # save x5~x31
19    .set n, 5
20    .rept 27
21        SAVE_GP %n
22        .set n, n+1
23    .endr
24    # we can use t0/t1/t2 freely, because they were saved on kernel stack
25    csrr t0, sstatus
26    csrr t1, sepc
27    sd t0, 32*8(sp)
28    sd t1, 33*8(sp)
29    # read user stack from sscratch and save it on the kernel stack
30    csrr t2, sscratch
31    sd t2, 2*8(sp)
32    # set input argument of trap_handler(cx: &mut TrapContext)
33    mv a0, sp # a0 point to trap context and as arguement of trap_handler
34    call trap_handler

```


__restore

- 先恢复CSR寄存器再恢复通用寄存器
- `csrrw sp, sscratch, sp` 此时sscratch设置为内核栈栈顶(由应用程序在执行前压入内核栈)

```
1  # os/src/trap/trap.S
2
3  .macro LOAD_GP n
4      ld x\n, \n*8(sp)
5  .endm
6
7  __restore:
8      # case1: start running app by __restore
9      # case2: back to U after handling trap
10     mv sp, a0
11     # now sp->kernel stack(after allocated), sscratch->user stack
12     # restore sstatus/sepc
13     ld t0, 32*8(sp)
14     ld t1, 33*8(sp)
15     ld t2, 2*8(sp)
16     csrw sstatus, t0
17     csrw sepc, t1
18     csrw sscratch, t2
19     # restore general-purpose registers except sp/tp
20     ld x1, 1*8(sp)
21     ld x3, 3*8(sp)
22     .set n, 5
23     .rept 27
24         LOAD_GP %n
25     .set n, n+1
26 .endr
27     # release TrapContext on kernel stack
28     addi sp, sp, 34*8
29     # now sp->kernel stack, sscratch->user stack
30     csrrw sp, sscratch, sp # 现在 sp 重新指向用户栈栈顶, sscratch 也依然保存进入
                               Trap 之前的状态并指向内核栈栈顶。
31     sret
```

trap_handler

```
1  // os/src/trap/mod.rs
2
```

```

3  #[no_mangle]
4  pub fn trap_handler(cx: &mut TrapContext) -> &mut TrapContext {
5      let scause = scause::read();
6      let stval = stval::read();
7      match scause.cause() {
8          Trap::Exception(Exception::UserEnvCall) => {
9              cx.sepc += 4;
10             cx.x[10] = syscall(cx.x[17], [cx.x[10], cx.x[11], cx.x[12]]) as
usize;
11         }
12         Trap::Exception(Exception::StoreFault) |
13         Trap::Exception(Exception::StorePageFault) => {
14             println!("[kernel] PageFault in application, kernel killed it.");
15             run_next_app();
16         }
17         Trap::Exception(Exception::IllegalInstruction) => {
18             println!("[kernel] IllegalInstruction in application, kernel
killed it.");
19             run_next_app();
20         }
21         _ => {
22             panic!("Unsupported trap {:?}, stval = {:#x}!", scause.cause(),
stval);
23         }
24     }
25     cx
26 }

```

Practice

扩展内核，能够统计多个应用的执行过程中系统调用编号和访问此系统调用的次数

- 增加一个SyscallNum的结构体记录发生系统调用的次数

```

1  // batch.rs
2  /**
3   * syscall num
4   */
5  pub struct SyscallNum {
6      num: [usize; SYSCALL_NUM],
7  }
8
9  impl SyscallNum {
10     /**

```

```

11     * get syscall num
12     */
13     pub fn get_syscall_num(&self, syscall_id: usize) -> usize {
14         self.num[syscall_id]
15     }
16     /**
17     * inc syscall num
18     */
19     pub fn inc_syscall_num(&mut self, syscall_id: usize) {
20         self.num[syscall_id] += 1;
21     }
22 }
23
24 lazy_static!{
25     // ch2 add begin
26     /**
27     * syscall use syscall num
28     */
29     pub static ref NUM: UPSafeCell<SyscallNum> = unsafe {
30         UPSafeCell::new({
31             SyscallNum {
32                 num: [0; SYSCALL_NUM],
33             }
34         })
35     };
36 }
37
38 // run_next_app
39 if current_app == APP_MANAGER.exclusive_access().num_app - 1 || current_app
== 0 {
40     println!(
41         "sys_write num: {}",
42         NUM.exclusive_access().get_syscall_num(0)
43     );
44     println!(
45         "sys_exit num: {}",
46         NUM.exclusive_access().get_syscall_num(1)
47     );
48 }

```

```

1 // syscall/mod.rs
2 use crate::batch::NUM;
3 /// handle syscall exception with `syscall_id` and other arguments
4 pub fn syscall(syscall_id: usize, args: [usize; 3]) -> isize {
5     match syscall_id {

```

```

6         SYSCALL_WRITE => {
7             let ret = sys_write(args[0], args[1] as *const u8, args[2]);
8             NUM.exclusive_access().inc_syscall_num(0);
9             ret
10        }
11        SYSCALL_EXIT => {
12            NUM.exclusive_access().inc_syscall_num(1);
13            sys_exit(args[0] as i32)
14        }
15        _ => panic!("Unsupported syscall_id: {}", syscall_id),
16    }

```

扩展内核，能够统计每个应用执行后的完成时间

- 在 `AppManager` 中增加应用执行时长字段；在 `run_next_app` 中获取系统时钟，记录应用时间

```

1  // batch.rs
2  /**
3   * AppManager
4   */
5  struct AppManager {
6      num_app: usize,
7      current_app: usize,
8      app_start: [usize; MAX_APP_NUM + 1],
9      // ch2 add begin
10     app_runtime: [u64; MAX_APP_NUM],
11     // ch2 add end
12 }
13
14 // run_next_app
15 // ch2 add begin
16
17 let time: u64;
18 unsafe {
19     asm!("rdtsc {0}", out(reg) time);
20 }
21 if current_app > 0 {
22     let runtime = time -
APP_MANAGER.exclusive_access().app_runtime[current_app - 1];
23     APP_MANAGER
24         .exclusive_access()
25         .set_app_runtime(current_app - 1, runtime);
26     println!(
27         "[kernel] app-{} runs {} cycles",
28         current_app - 1,

```

```

29     APP_MANAGER
30     .exclusive_access()
31     .get_app_runtime(current_app - 1)
32 );
33 }

```

sys_write 仅能输出位于程序本身内存空间内的数据，否则报错

- 构建ReliableAddr结构体记录应用程序起始和终止地址

```

1  // batch.rs
2  /**
3   * Reliable address of program
4   */
5  pub struct ReliableAddr {
6      start: usize,
7      end: usize,
8  }
9
10 impl ReliableAddr {
11     /**
12      * get reliable start addr
13      */
14     pub fn get_reliable_start(&self) -> usize {
15         self.start
16     }
17     /**
18      * get reliable end addr
19      */
20     pub fn get_reliable_end(&self) -> usize {
21         self.end
22     }
23 }
24
25 lazy_static!{
26     /**
27      * reliable addr init
28      */
29     pub static ref RE_ADDR:UPSafeCell<ReliableAddr> = unsafe {
30         UPSafeCell::new({
31             ReliableAddr {
32                 start: APP_BASE_ADDRESS,
33                 end: APP_BASE_ADDRESS + APP_SIZE_LIMIT,
34             }
35         })

```

```

36     };
37 }
38
39 // run_next_app
40 // 注意数据需要使用clone()
41 let mut reliable_addr = RE_ADDR.exclusive_access();
42 reliable_addr.start = app_manager.app_start[current_app].clone();
43 reliable_addr.end = app_manager.app_start[current_app + 1].clone();
44 drop(reliable_addr);

```

```

1 // syscall/mod.rs
2 use crate::batch::RE_ADDR;
3
4 /// handle syscall exception with `syscall_id` and other arguments
5 pub fn syscall(syscall_id: usize, args: [usize; 3]) -> isize {
6     match syscall_id {
7         SYSCALL_WRITE => {
8             let start_addr = RE_ADDR.exclusive_access().get_reliable_start();
9             let end_addr = RE_ADDR.exclusive_access().get_reliable_end();
10            if args[1] < start_addr {
11                return -1;
12            }
13            if args[1] >= end_addr {
14                return -1;
15            }
16            if args[1] + args[2] >= end_addr {
17                return -1;
18            }
19            let ret = sys_write(args[0], args[1] as *const u8, args[2]);
20            NUM.exclusive_access().inc_syscall_num(0);
21            ret
22        }
23        SYSCALL_EXIT => {
24            NUM.exclusive_access().inc_syscall_num(1);
25            sys_exit(args[0] as i32)
26        }
27        _ => panic!("Unsupported syscall_id: {}", syscall_id),
28    }
29 }

```