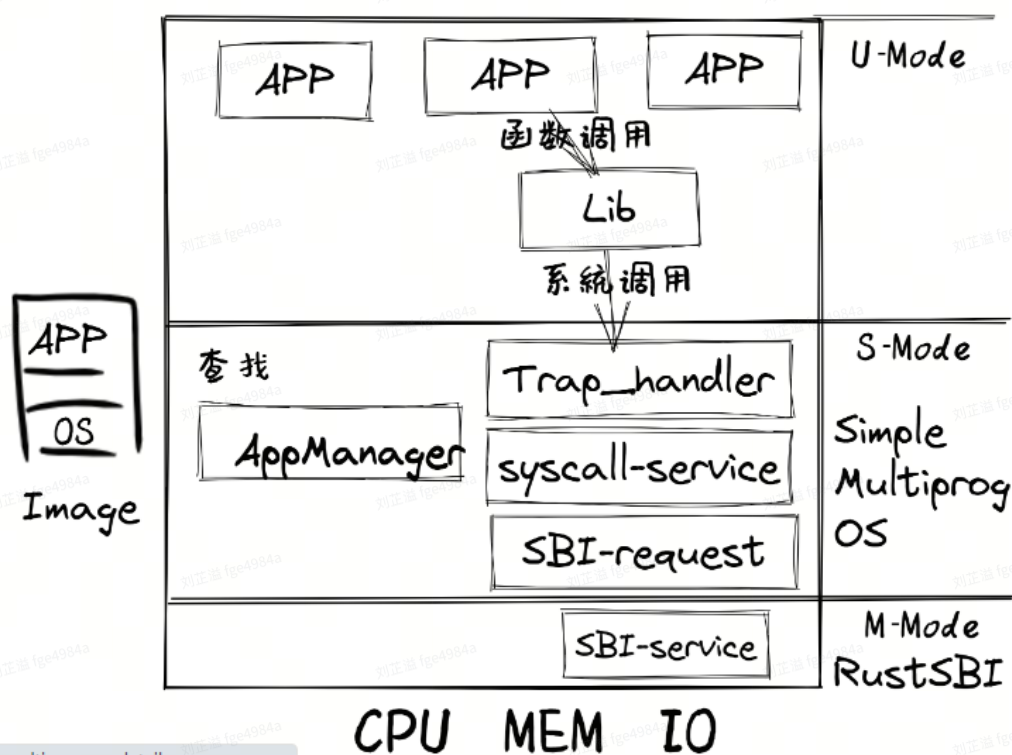


# TimeSharingOS

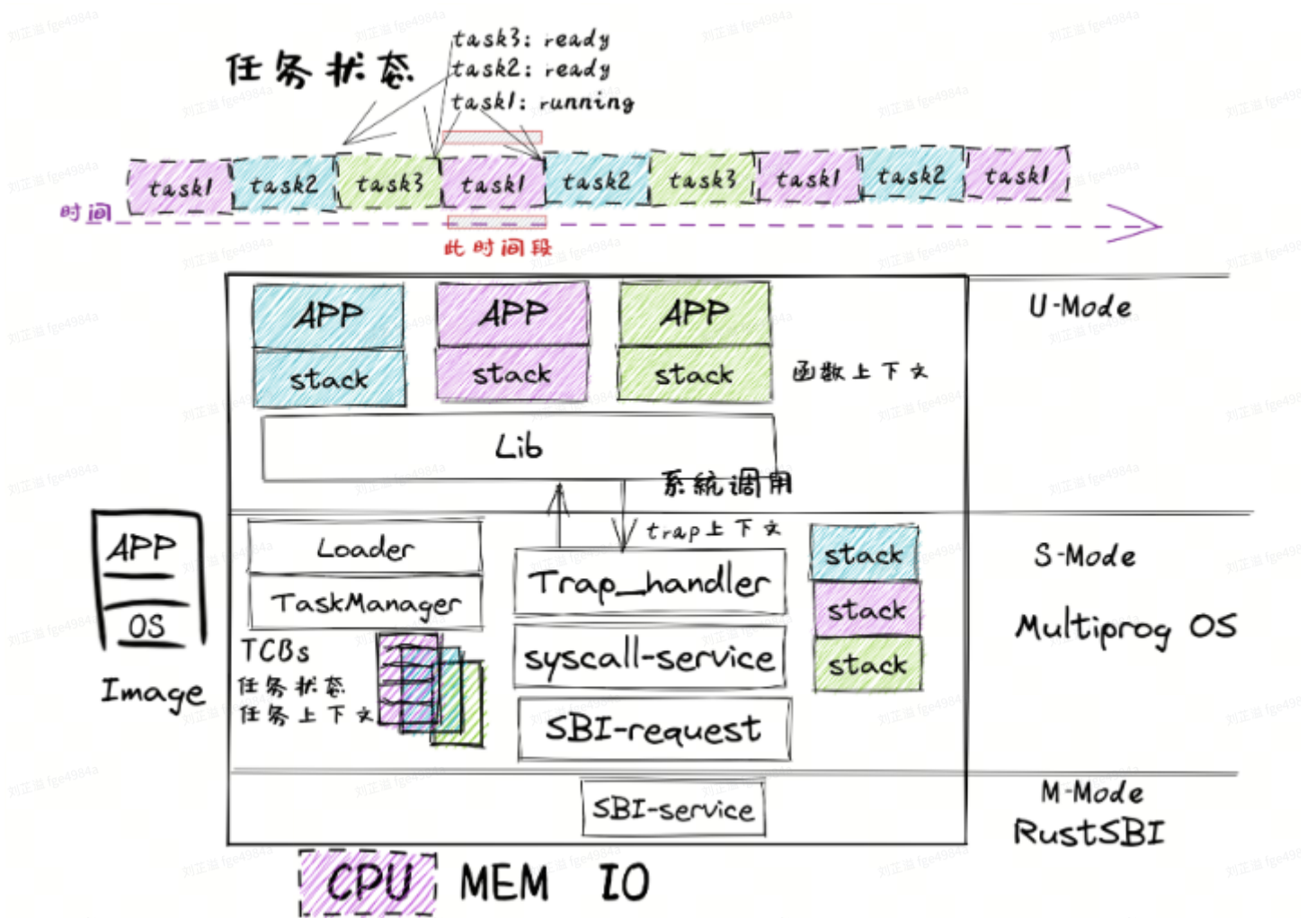
## TimesharingOS

## MultiprogOS



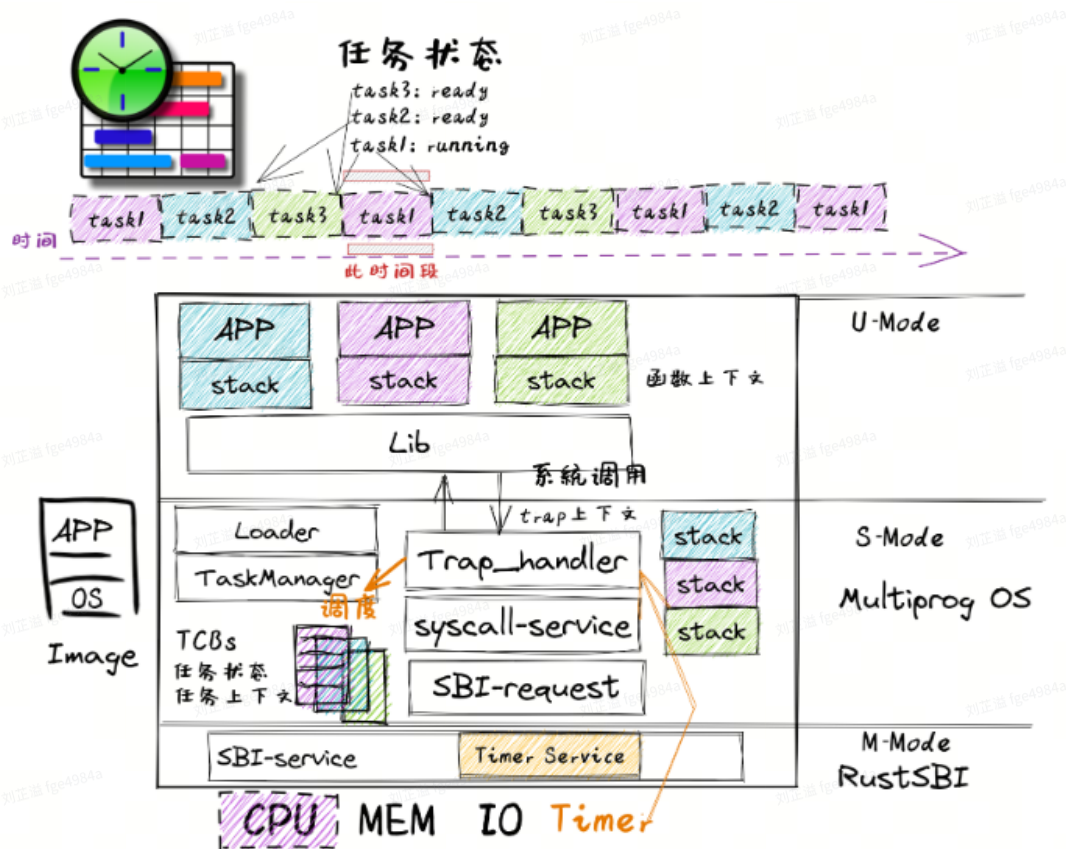
- Qemu把包含多个app的列表和MultiprogOS的image镜像加载到内存中，RustSBI（bootloader）完成基本的硬件初始化后，跳转到MultiprogOS起始位置
- MultiprogOS首先进行正常运行前的初始化工作，即建立栈空间和清零bss段，然后通过改进的AppManager内核模块从app列表中把**所有app都加载到内存中**，并按指定顺序让app在用户态一个接一个地执行。
- app在执行过程中，会通过系统调用的方式得到MultiprogOS提供的OS服务，如输出字符串等。

## CoopOS



- CoopOS进一步改进了AppManager内核模块，把它拆分为负责**加载应用的 Loader 内核模块**和管理应用运行过程的TaskManager内核模块。
- TaskManager通过task任务控制块来管理应用程序的执行过程，**支持应用程序主动放弃CPU并切换到另一个应用继续执行**，从而提高系统整体执行效率。
- 应用程序在运行时有自己所在的内存空间和栈，确保被切换时相关信息不会被其他应用破坏。如果当前应用程序正在运行，则该应用对应的任务处于运行（Running）状态；如果该应用主动放弃处理器，则该应用对应的任务处于就绪（Ready）状态。
- 操作系统进行任务切换时，需要把要暂停任务的上下文（即任务用到的通用寄存器）保存起来，把要继续执行的任务的上下文恢复为暂停前的内容，这样就能让不同的应用协同使用处理器了。

## TimersharingOS



- TimesharingOS最大的变化是改进了 Trap\_handler 内核模块，**支持时钟中断**，从而可以抢占应用的执行。
- 并通过进一步改进 TaskManager 内核模块，提供**任务调度功能**，这样可以在收到时钟中断后统计任务的使用时间片，如果任务的时间片用完后，则切换任务。从而可以公平和高效地分时执行多个应用，提高系统的整体效率。

## 多道程序放置与加载

### 放置

- 每个app的起始地址为 `hex(base_address+step*app_id)`，实现起始地址的不同

```

1  # user/build.py
2
3  import os
4
5  base_address = 0x80400000
6  step = 0x20000
7  linker = 'src/linker.ld'
8
9  app_id = 0
10 apps = os.listdir('src/bin')
11 apps.sort()
12 for app in apps:
13     app = app[:app.find('.')]

```

```

14     lines = []
15     lines_before = []
16     with open(linker, 'r') as f:
17         for line in f.readlines():
18             lines_before.append(line)
19             line = line.replace(hex(base_address),
20                                 hex(base_address+step*app_id))
21             lines.append(line)
22     with open(linker, 'w+') as f:
23         f.writelines(lines)
24     os.system('cargo build --bin %s --release' % app)
25     print('[build.py] application %s start with address %s' %(app,
26                                     hex(base_address+step*app_id)))
27     with open(linker, 'w+') as f:
28         f.writelines(lines_before)
29     app_id = app_id + 1

```

## 加载

- 第*i*个应用加载到不同的物理地址上  $\text{APP\_BASE\_ADDRESS} + i * \text{APP\_SIZE\_LIMIT}$

```

1 // os/src/loader.rs
2
3 fn get_base_i(app_id: usize) -> usize {
4     APP_BASE_ADDRESS + app_id * APP_SIZE_LIMIT
5 }
6
7 pub fn load_apps() {
8     extern "C" { fn _num_app(); }
9     let num_app_ptr = _num_app as usize as *const usize;
10    let num_app = get_num_app();
11    let app_start = unsafe {
12        core::slice::from_raw_parts(num_app_ptr.add(1), num_app + 1)
13    };
14    // load apps
15    for i in 0..num_app {
16        let base_i = get_base_i(i);
17        // clear region
18        (base_i..base_i + APP_SIZE_LIMIT).for_each(|addr| unsafe {
19            (addr as *mut u8).write_volatile(0)
20        });
21        // load app from data section to memory
22        let src = unsafe {
23            core::slice::from_raw_parts(
24                app_start[i] as *const u8,
25                app_start[i + 1] - app_start[i]

```

```

26         )
27     };
28     let dst = unsafe {
29         core::slice::from_raw_parts_mut(base_i as *mut u8, src.len())
30     };
31     dst.copy_from_slice(src);
32 }
33 unsafe {
34     asm!("fence.i");
35 }
36 }

```

## 任务切换

- 在内核中这种机制是在 `__switch` 函数中实现的。任务切换支持的场景是：一个应用在运行途中便会主动或被动交出 CPU 的使用权，此时它只能暂停执行，等到内核重新给它分配处理器资源之后才能恢复并继续执行。
- 任务切换是来自**两个不同应用在内核中的 Trap 控制流**之间的切换。当一个应用 Trap 到 S 模式的操作系统内核中进行进一步处理（即进入了操作系统的 Trap 控制流）的时候，其 Trap 控制流可以调用一个特殊的 `__switch` 函数。
- 这个函数表面上就是一个普通的函数调用：在 `__switch` 返回之后，将继续从调用该函数的位置继续向下执行。但是其间却隐藏着复杂的控制流切换过程。具体来说，调用 `__switch` 之后直到它返回前的这段时间，
  - 原 Trap 控制流 A 会先被暂停并被切换出去，CPU 转而运行另一个应用在内核中的 Trap 控制流 B。
  - 然后在某个合适的时机，原 Trap 控制流 A 才会从某一条 Trap 控制流 C（很有可能不是它之前切换到的 B）切换回来继续执行并最终返回。
  - 不过，从实现的角度讲，`__switch` 函数和一个普通的函数之间的核心差别仅仅是它会 **换栈**。

## 设计与实现

- 对于当前正在执行的任务的 Trap 控制流，我们用一个名为 `current_task_cx_ptr` 的变量来保存放置当前任务上下文的地址；而用 `next_task_cx_ptr` 的变量来保存放置下一个要执行任务的上下文的地址。利用 C 语言的引用来描述的话就是：

```

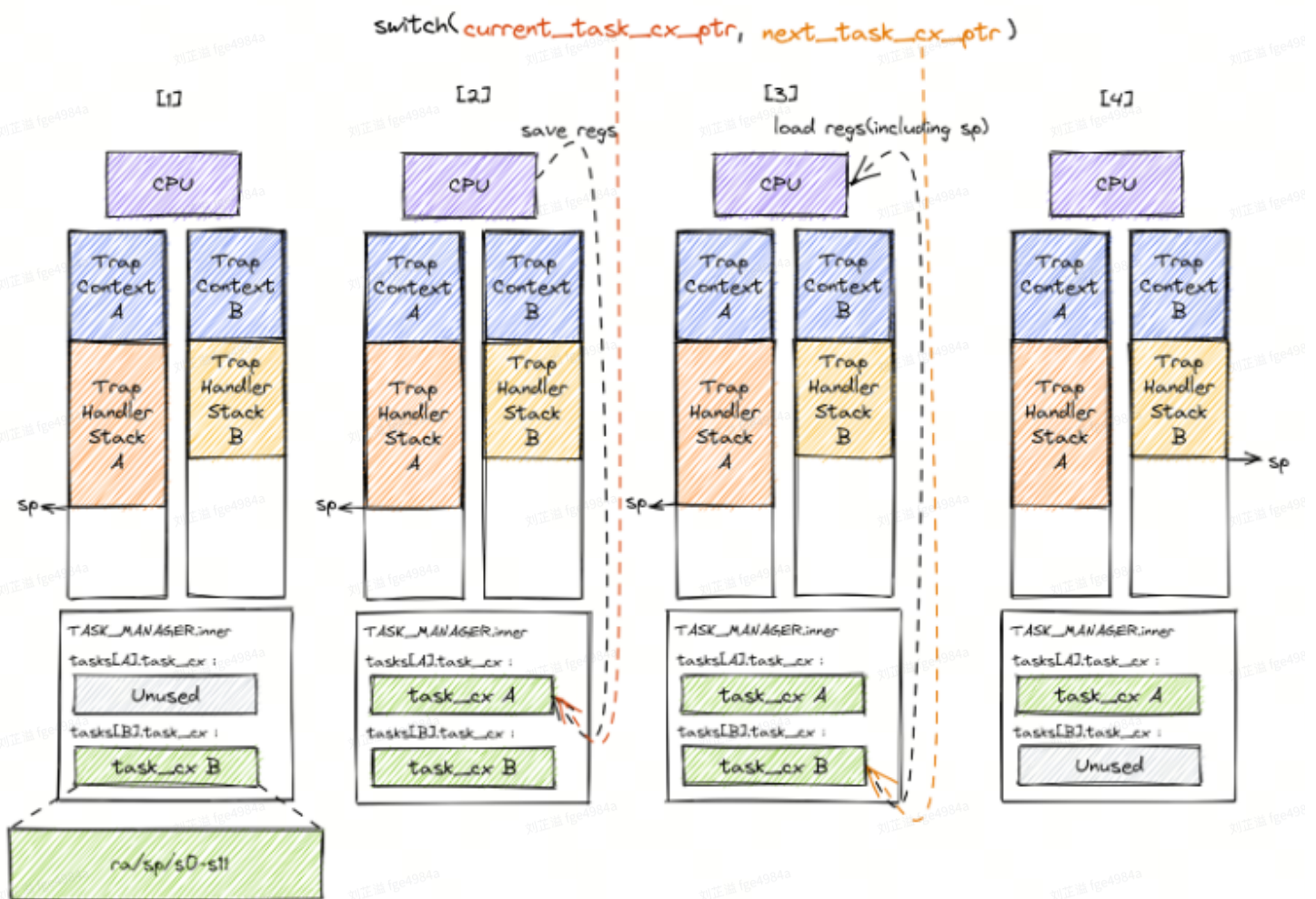
1 TaskContext *current_task_cx_ptr = &tasks[current].task_cx;
2 TaskContext *next_task_cx_ptr    = &tasks[next].task_cx;

```

- 任务切换包含4个阶段



- a. 在 Trap 控制流 A 调用 `__switch` 之前，A 的内核栈上只有 Trap 上下文和 Trap 处理函数的调用栈信息，而 B 是之前被切换出去的；
  - b. A 在 A 任务上下文空间在里面保存 CPU 当前的寄存器快照；
  - c. 读取 `next_task_cx_ptr` 指向的 B 任务上下文，根据 B 任务上下文保存的内容来恢复 `ra` 寄存器、`s0~s11` 寄存器以及 `sp` 寄存器。只有这一步做完后，`__switch` 才能做到一个函数跨两条控制流执行，即通过换栈也就实现了控制流的切换。
  - d. 上一步寄存器恢复完成后，可以看到通过恢复 `sp` 寄存器换到了任务 B 的内核栈上，进而实现了控制流的切换。这就是为什么 `__switch` 能做到一个函数跨两条控制流执行。此后，当 CPU 执行 `ret` 汇编伪指令完成 `__switch` 函数返回后，任务 B 可以从调用 `__switch` 的位置继续向下执行。
- 从结果来看，我们看到 A 控制流和 B 控制流的状态发生了互换，A 在保存任务上下文之后进入暂停状态，而 B 则恢复了上下文并在 CPU 上继续执行。



## \_\_switch实现

```

1 // os/src/task/switch.rs
2
3 global_asm!(include_str!("switch.S"));
4

```

```

5  use super::TaskContext;
6
7  extern "C" {
8      pub fn __switch(
9          current_task_cx_ptr: *mut TaskContext,
10         next_task_cx_ptr: *const TaskContext
11     );
12 }

```

```

1  # os/src/task/switch.S
2
3  .altmacro
4  .macro SAVE_SN n
5      sd s\n, (\n+2)*8(a0)
6  .endm
7  .macro LOAD_SN n
8      ld s\n, (\n+2)*8(a1)
9  .endm
10  .section .text
11  .globl __switch
12  __switch:
13      # 阶段 [1]
14      # __switch(
15      #     current_task_cx_ptr: *mut TaskContext,
16      #     next_task_cx_ptr: *const TaskContext
17      # )
18      # 阶段 [2]
19      # save kernel stack of current task
20      sd sp, 8(a0)
21      # save ra & s0~s11 of current execution
22      sd ra, 0(a0)
23      .set n, 0
24      .rept 12
25          SAVE_SN %n
26          .set n, n + 1
27      .endr
28      # 阶段 [3]
29      # restore ra & s0~s11 of next execution
30      ld ra, 0(a1)
31      .set n, 0
32      .rept 12
33          LOAD_SN %n
34          .set n, n + 1
35      .endr
36      # restore kernel stack of next task

```

```

37     ld sp, 8(a1)
38     # 阶段 [4]
39     ret

```

## TaskContext

```

1  // os/src/task/context.rs
2
3  pub struct TaskContext {
4      ra: usize,
5      sp: usize,
6      s: [usize; 12],
7  }

```

## 任务管理器

### 任务控制块(TCB)

```

1  // os/src/task/task.rs
2
3  #[derive(Copy, Clone)]
4  pub struct TaskControlBlock {
5      pub task_status: TaskStatus,
6      pub task_cx: TaskContext,
7  }

```

- `TaskContext::goto_restore` 将 `ra` 设置为 `__restore`；因为切换任务时我们设定内核栈栈中包含 `__alltraps` 和 `trap_handler` 的栈帧
- 当执行第一个程序时，我们需要向内核栈中压入一个 `TrapContext`

```

1  // os/src/loader.rs
2
3  pub fn init_app_cx(app_id: usize) -> usize {
4      KERNEL_STACK[app_id].push_context(
5          TrapContext::app_init_context(get_base_i(app_id),
6          USER_STACK[app_id].get_sp()),
7      )
8  }
9  // os/src/task/mod.rs

```



```

10 pub struct TaskManager {
11     num_app: usize,
12     inner: UPSafeCell<TaskManagerInner>,
13 }
14
15 struct TaskManagerInner {
16     tasks: [TaskControlBlock; MAX_APP_NUM],
17     current_task: usize,
18 }
19
20 lazy_static! {
21     pub static ref TASK_MANAGER: TaskManager = {
22         let num_app = get_num_app();
23         let mut tasks = [
24             TaskControlBlock {
25                 task_cx: TaskContext::zero_init(),
26                 task_status: TaskStatus::UnInit
27             };
28             MAX_APP_NUM
29         ];
30         for i in 0..num_app {
31             tasks[i].task_cx = TaskContext::goto_restore(init_app_cx(i));
32             tasks[i].task_status = TaskStatus::Ready;
33         }
34         TaskManager {
35             num_app,
36             inner: unsafe { UPSafeCell::new(TaskManagerInner {
37                 tasks,
38                 current_task: 0,
39             }) },
40         }
41     };
42 }

```

## sys\_yield和sys\_exit系统调用

```

1 // os/src/syscall/process.rs
2
3 use crate::task::suspend_current_and_run_next;
4
5 pub fn sys_yield() -> isize {
6     suspend_current_and_run_next();
7     0
8 }
9 // os/src/syscall/process.rs

```

```

10
11 use crate::task::exit_current_and_run_next;
12
13 pub fn sys_exit(exit_code: i32) -> ! {
14     println!("[kernel] Application exited with code {}", exit_code);
15     exit_current_and_run_next();
16     panic!("Unreachable in sys_exit!");
17 }

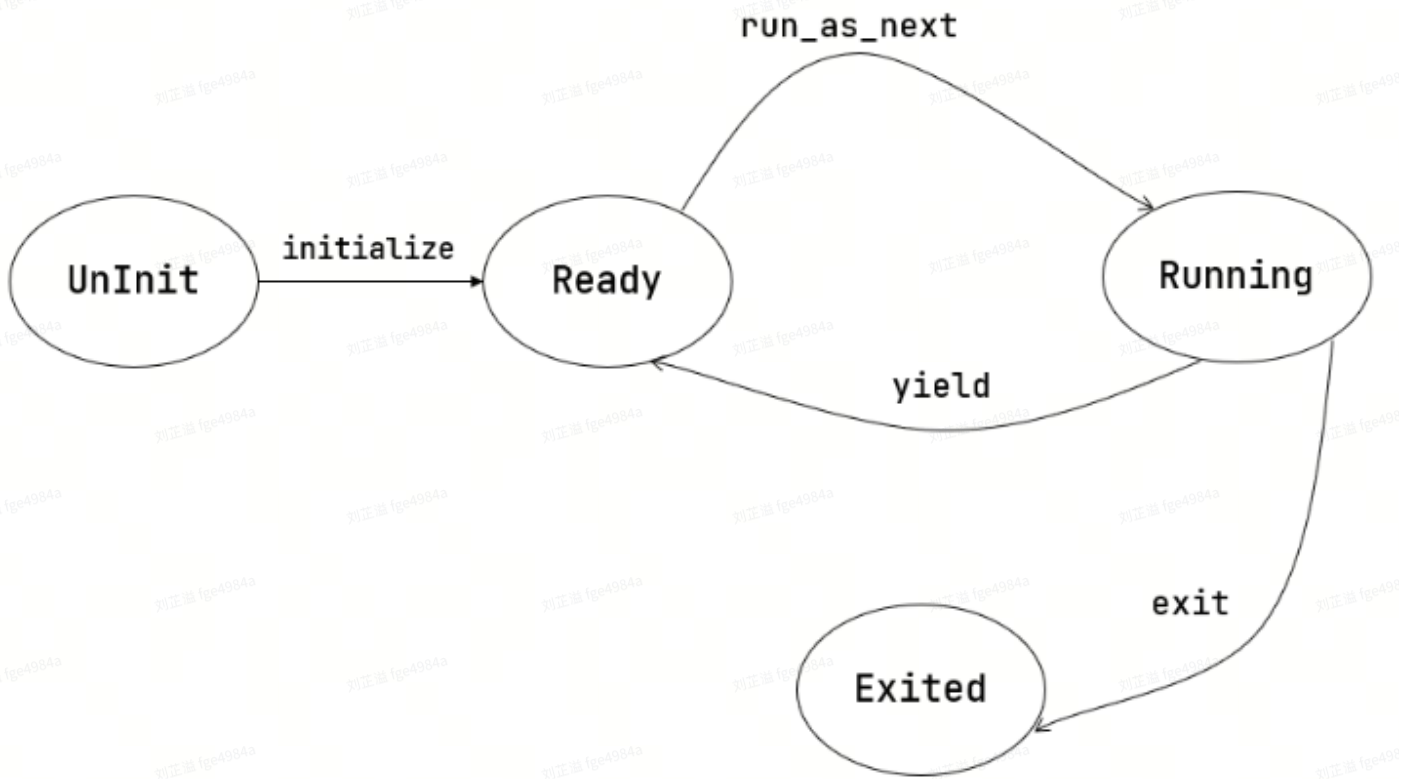
```

- `suspend_current_and_run_next` 和 `exit_current_and_run_next` 均是切换当前 Task 的运行状态，切换到下一个应用

```

1 // os/src/task/mod.rs
2
3 fn mark_current_suspended() {
4     TASK_MANAGER.mark_current_suspended();
5 }
6
7 fn mark_current_exited() {
8     TASK_MANAGER.mark_current_exited();
9 }
10
11 impl TaskManager {
12     fn mark_current_suspended(&self) {
13         let mut inner = self.inner.borrow_mut();
14         let current = inner.current_task;
15         inner.tasks[current].task_status = TaskStatus::Ready;
16     }
17
18     fn mark_current_exited(&self) {
19         let mut inner = self.inner.borrow_mut();
20         let current = inner.current_task;
21         inner.tasks[current].task_status = TaskStatus::Exited;
22     }
23 }

```



## run\_next\_task

```
1 // os/src/task/mod.rs
2
3 fn run_next_task() {
4     TASK_MANAGER.run_next_task();
5 }
6
7 impl TaskManager {
8     fn run_next_task(&self) {
9         if let Some(next) = self.find_next_task() {
10             let mut inner = self.inner.exclusive_access();
11             let current = inner.current_task;
12             inner.tasks[next].task_status = TaskStatus::Running;
13             inner.current_task = next;
14             let current_task_cx_ptr = &mut inner.tasks[current].task_cx as
15             *mut TaskContext;
16             let next_task_cx_ptr = &inner.tasks[next].task_cx as *const
17             TaskContext;
18             drop(inner);
19             // before this, we should drop local variables that must be
20             // dropped manually
21             unsafe {
22                 __switch(
23                     current_task_cx_ptr,
24                     next_task_cx_ptr,
```

```

22         );
23     }
24     // go back to user mode
25 } else {
26     panic!("All applications completed!");
27 }
28 }
29 }

```

## 时间片轮转调度(Round-Robin)

### 时钟中断与计时器

- 在 RISC-V 64 架构上，该计数器保存在一个 64 位的 CSR `mtime` 中，我们无需担心它的溢出问题，在内核运行全程可以认为它是一直递增的。
- 另外一个 64 位的 CSR `mtimecmp` 的作用是：一旦计数器 `mtime` 的值超过了 `mtimecmp`，就会触发一次时钟中断。这使得我们可以方便的通过设置 `mtimecmp` 的值来决定下一次时钟中断何时触发。
- 可惜的是，它们都是 M 特权级的 CSR，而我们的内核处在 S 特权级，是不被允许直接访问它们的。好在运行在 M 特权级的 SEE（这里是 RustSBI）已经预留了相应的接口，我们可以调用它们来间接实现计时器的控制
- 常数 `CLOCK_FREQ` 是一个预先获取到的各平台不同的时钟频率，单位为赫兹，也就是一秒钟之内计数器的增量。

```

1 // os/src/timer.rs
2
3 use riscv::register::time;
4
5 pub fn get_time() -> usize {
6     time::read()
7 }
8 // os/src/sbi.rs
9
10 const SBI_SET_TIMER: usize = 0;
11
12 pub fn set_timer(timer: usize) {
13     sbi_call(SBI_SET_TIMER, timer, 0, 0);
14 }

```

- `set_next_trigger` 设置下一个时钟中断
- `get_time_us` 以微秒为单位返回当前计数器的值

```

1 // os/src/timer.rs
2
3 use crate::config::CLOCK_FREQ;
4 const TICKS_PER_SEC: usize = 100;
5
6 pub fn set_next_trigger() {
7     set_timer(get_time() + CLOCK_FREQ / TICKS_PER_SEC);
8 }
9
10 // os/src/timer.rs
11
12 const MICRO_PER_SEC: usize = 1_000_000;
13
14 pub fn get_time_us() -> usize {
15     time::read() / (CLOCK_FREQ / MICRO_PER_SEC)
16 }

```

## 抢占式调度

```

1 // os/src/trap/mod.rs
2
3 match scause.cause() {
4     Trap::Interrupt(Interrupt::SupervisorTimer) => {
5         set_next_trigger();
6         suspend_current_and_run_next();
7     }
8 }

```

- 我们只需在 `trap_handler` 函数下新增一个条件分支跳转，当发现触发了一个 S 特权级时钟中断的时候，首先重新设置一个 10ms 的计时器，然后调用上一小节提到的 `suspend_current_and_run_next` 函数暂停当前应用并切换到下一个。
- 初始化设置

```

1 // os/src/main.rs
2
3 #[no_mangle]
4 pub fn rust_main() -> ! {
5     clear_bss();
6     println!("[kernel] Hello, world!");
7     trap::init();
8     loader::load_apps();
9 }

```



```

9
10     trap::enable_timer_interrupt();
11     timer::set_next_trigger();
12
13     task::run_first_task();
14     panic!("Unreachable in rust_main!");
15 }
16
17 // os/src/trap/mod.rs
18
19 use riscv::register::sie;
20
21 pub fn enable_timer_interrupt() {
22     unsafe { sie::set_stimer(); }
23 }

```

## sleep

- 目前在等待某些事件的时候仍然需要 `yield`，其中一个原因是为了节约 CPU 计算资源，另一个原因是当事件依赖于其他的应用的时候，由于只有一个 CPU，当前应用的等待可能永远不会结束。这种情况下需要先将它切换出去，使得其他的应用到达它所期待的状态并满足事件的生成条件，再切换回来。
- 这里我们先通过 `yield` 来优化 **轮询** (Busy Loop) 过程带来的 CPU 资源浪费。在 `03sleep` 这个应用中：

```

1 // user/src/bin/03sleep.rs
2
3 #[no_mangle]
4 fn main() -> i32 {
5     let current_timer = get_time();
6     let wait_for = current_timer + 3000;
7     while get_time() < wait_for {
8         yield_();
9     }
10    println!("Test sleep OK!");
11    0
12 }

```