

# 操作系统内核设计

Del0n1x

哈尔滨工业大学(深圳)

2025 年 6 月

## 目录

第 1 章 概述 .....	1
1.1 项目介绍 .....	1
第 2 章 进程管理 .....	2
2.1 概述 .....	2
2.2 任务调度 .....	2
2.2.1 异步并不是"银弹" .....	3
2.3 任务调度队列与执行器 .....	3
2.4 多核心 CPU 管理 .....	5
2.5 任务控制块 .....	6
2.5.1 进程和线程联系 .....	7
2.5.2 任务的状态 .....	8
第 3 章 中断与异常处理 .....	9
3.1 特权级切换 .....	9
3.2 处理过程 .....	10
3.2.1 内核中断异常处理 .....	10
3.2.2 用户中断异常处理 .....	11
3.2.3 返回用户态 .....	12
第 4 章 内存管理 .....	13
4.1 物理内存管理 .....	13
4.1.1 内核动态内存分配器 .....	13
4.1.2 物理页分配器 .....	14
4.2 地址空间 .....	14
4.2.1 地址空间布局 .....	14
4.2.2 Boot 阶段高位映射 .....	15
4.2.3 地址空间管理 .....	16
4.3 缺页异常处理 .....	17
4.3.1 CoW 写时复制技术 .....	17
4.3.2 懒分配技术 .....	18
4.3.3 用户地址检查与零拷贝技术 .....	19

第 5 章 文件系统 .....	20
5.1 虚拟文件系统 .....	20
5.1.1 SuperBlock .....	21
5.1.2 Inode .....	22
5.1.3 Dentry .....	23
5.1.4 File .....	26
5.2 磁盘文件系统 .....	26
5.2.1 EXT4 文件系统 .....	26
5.3 非磁盘文件系统 .....	27
5.3.1 procfs .....	27
5.3.2 devfs .....	27
5.4 页缓存 .....	27
5.5 其他数据结构 .....	28
5.5.1 FdTable .....	28
第 6 章 进程间通信 .....	29
6.1 信号机制 .....	29
6.2 信号传输 .....	30
6.3 信号处理 .....	31
第 7 章 时钟模块 .....	32
7.1 定时器队列 .....	32
7.1.1 时间轮设计 .....	33
7.2 定时器 .....	34
第 8 章 网络模块 .....	35
8.1 Socket 套接字 .....	35
8.2 Ethernet 设备 .....	37
8.3 传输层——UDP 与 TCP .....	38
8.4 Port 端口分配 .....	39
第 9 章 设备 .....	40
9.1 设备树 .....	41
第 10 章 硬件抽象层 .....	42
10.1 硬件抽象层 (Hardware Abstraction Layer, HAL) 总览 .....	42
10.2 处理器访问接口 .....	42

10.3	内核入口例程 .....	43
10.4	内存管理单元与地址空间 .....	43
10.4.1	物理内存 .....	43
10.4.2	分页地址翻译模式 .....	43
10.4.3	页表 .....	44
10.4.4	直接映射窗口 .....	44
10.4.5	TLB 重填 .....	45
第 11 章	总结与展望 .....	45
11.1	工作总结 .....	45
11.2	未来计划 .....	45
11.3	参考 .....	46

## 第 1 章 概述

### 1.1 项目介绍

Del0n1x 是一个使用 Rust 语言编写的同时适配 RISC-V64 和 LoongArch64 的宏内核，支持多核运行与无栈协程进程调度，拥有完善的内存管理和信号传递机制。在软硬件交换层面，我们实现了自己的 HAL 层，统一调用接口，能够同时支持 RISC-V64 和 LoongArch64 指令架构。

Del0n1x 致力于实现高效清晰的代码逻辑，遵守 System Manual 手册，实现 105 个相关系统调用，并且对于其中大部分系统调用做了相对完善的错误检查和处理机制，这为我们后来适配 ltp 带来了便利。

Del0n1x 初赛阶段的内核主要模块和完成情况如下表格：

模块	完成情况
HAL 模块	实现自己的 HAL 代码库，支持 riscv64 和 loongarch64 双架构
进程管理	无栈协程调度，支持全局统一的 executor 调度器 实现多线程的资源回收 统一进程和线程的数据结构
文件系统	实现 dentry 构建目录树 实现页缓存和 dentry 缓存加快读写
内存管理	实现 COW、懒加载内存优化 实现地址空间共享
时钟模块	实现时间轮混合最小堆的数据结构管理方式 支持定时器唤醒机制
信号系统	支持处理用户自定义信号和 sigreturn 机制
多核管理	支持多核运行，但是还有部分 bug
网络模块	初步完成网络模块相关代码，由于时间原因还没有适配通过网络测例

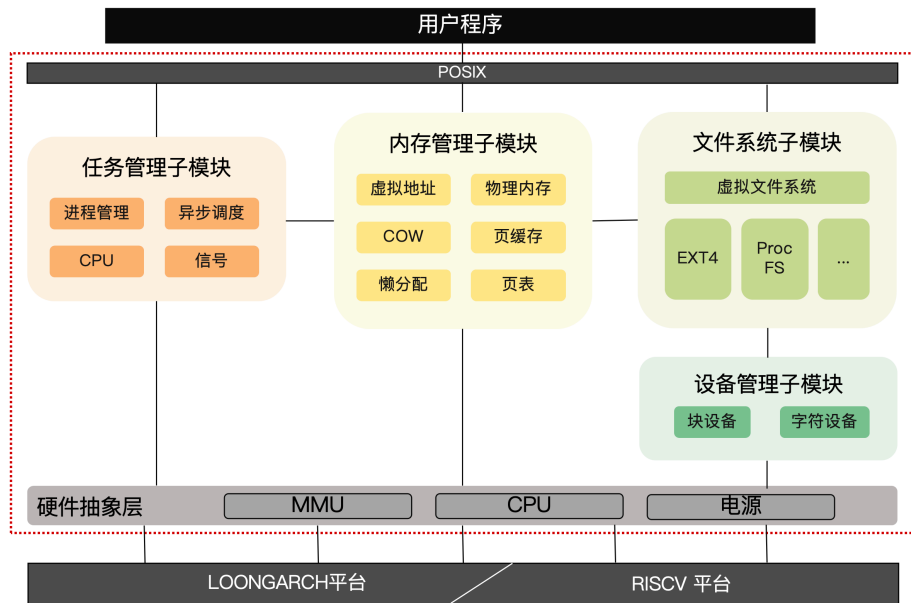


图 1-1 Del0n1x 整体架构图

## 第 2 章 进程管理

### 2.1 概述

Del0n1x 操作系统采用无栈协程作为核心任务管理模型，该设计基于用户级轻量级线程理念，允许任务在执行过程中挂起并恢复。无栈协程与有栈协程的核心差异体现在上下文管理机制上：在 rcore 和 xv6 中，采用有栈方式进行任务切换，每次任务切换需要调用 `__switch` 函数，每个任务有独立栈空间并在切换时保存完整栈帧及 `s0-s11`、`ra`、`sp` 等 14 个 RISC-V 寄存器，而 Del0n1x 的无栈协程不依赖独立栈结构，转通过状态机管理上下文状态，任务挂起时仅需保存当前执行位置和关键状态变量，显著降低内存开销与切换延迟。

在传统操作系统中，进程调度需切换用户级上下文、寄存器上下文和系统级上下文（含内核栈与页表），这种完整上下文切换必须通过内核态系统调用完成。但是 Del0n1x 采用共享内核栈架构，所有任务复用同一内核栈空间，任务切换时不涉及内核栈切换，这极大地降低了任务切换的开销。

### 2.2 任务调度

在无栈协程中，任务调度并不需要在栈上保存任务栈帧，而是将必要的中间信息保存为状态机，放置在堆空间进行保存，这样不仅减少了任务调度开销，同时提高了调度过程中的安全性，降低栈溢出的风险。

在 Del0n1x 中，当对异步函数调用 `.await` 方法时，`async-task` 库会首次调用 `Future` 的 `poll` 方法，如果 `poll` 返回的结果是 `Pending`，那么该任务将被挂起，`await` 将控制权交给调度器，以便另一个任务可以继续进行。任务调度如下图：

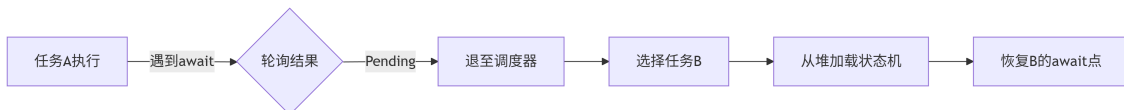


图 2-1 任务调度

### 2.2.1 异步并不是“银弹”

我们在第一次接触异步调度时觉得，既然异步调度效率高，那为何不把所有的系统调用设计为异步形式呢？后来在深入了解到异步机制后发现，我们陷入了一个常见误区：把异步当成万能银弹。

首先异步的核心价值在于解决 I/O 瓶颈。I/O 的延迟远高于 CPU 计算，磁盘/网络操作耗时可能是微秒（ $\mu s$ ）甚至毫秒（ $ms$ ）级，而 CPU 指令是纳秒（ $ns$ ）级。异步 I/O 允许 CPU 在等待慢速 I/O 时执行其他任务。CPU 密集型操作无等待需求，如 `getpid()`、`sched_yield()` 等系统调用本身不阻塞，异步化不会带来收益，反而增加调度开销。若一个操作本身能在短时间内完成，同步调用直接返回结果的开销远小于异步回调的上下文切换。

其次状态管理困难。异步操作需通过回调和细致的 `Future` 设计，如果想要精细的控制 `poll` 中代码处理流程，需要自己实现 `Future` 的 `poll` 轮循进行状态管理，这无疑会复杂化内核实现与 bug 调试。

## 2.3 任务调度队列与执行器

在任务调度队列实现中，Del0n1x 将调度队列分为 FIFO 和 PRIO 队列，感谢优秀作品 `Plntry` 对 `async-task` 库做出的 `Pr`，使得在任务调度过程中可以获取调度信息，通过调度信息，我们可以对任务做出更加精细的控制。如果任务在运行时被唤醒，则将其加入 FIFO 队列，其他的就放入 PRIO 队列进行管理。

```
1 struct TaskQueue {
2     normal: SpinNoIrqLock<VecDeque<Runnable>>,
3     prior: SpinNoIrqLock<VecDeque<Runnable>>,
4 }
5
6 ....
7 // 任务入队逻辑
8 if info.woken_while_running {
9     queue.push_normal(runnable);
10 } else {
11     queue.push_prior(runnable);
12 }
13
14 ....
```

代码 2-1 任务队列结构

在 Del0n1x 中，我们使用统一的 TaskFuture 封装了任务。对于用户任务，在 poll 轮循中实现了任务的切换调度。当任务 checkin 时，需要在修改 TCB 的时间戳记录调度时间，然后切换 CPU 中运行任务和切换页表。当任务 checkout 时，需要判断浮点寄存器状态是否为 dirty 以确定是否保存浮点寄存器，然后清空 CPU 当前任务，并记录任务 checkout 时间。对于内核任务，Del0n1x 并没有设计任务切换，而是让该任务一直 poll，直到任务结束，这类任务主要是 shell 程序。



```

1  pub enum TaskFuture<F: Future<Output = ()> + Send + 'static> {
2      UserTaskFuture {
3          task: Arc<TaskControlBlock>,
4          future: F,
5      },
6      KernelTaskFuture {
7          future: F,
8      },
9  }
10
11 fn poll(
12     self: Pin<&mut Self>,
13     cx: &mut core::task::Context<'_,>,
14 ) -> core::task::Poll<Self::Output> {
15     let this = unsafe { self.get_unchecked_mut() };
16
17     match this {
18         TaskFuture::UserTaskFuture { task, future } => {
19             let processor = get_current_cpu();
20             processor.user_task_checkin(task); // 用户任务 checkin
21             let ret = unsafe { Pin::new_unchecked(future).poll(cx) };
22             processor.user_task_checkout(task); // 用户任务 checkout
23             ret
24         }
25         TaskFuture::KernelTaskFuture { future } => {
26             unsafe { Pin::new_unchecked(future).poll(cx) }
27         }
28     }
29 }

```

代码 2-2 任务 Future 结构与 poll 实现

spawn\_user\_task 可以设置一个用户任务。Del0n1x 将用户任务的 future 设置为 trap\_loop 循环，负责处理任务在用户态和内核态之间的切换，直到任务结束。执行 executor::spawn(future) 将任务挂入全局队列中等待被调度。

```

1  pub fn spawn_user_task(user_task: Arc<TaskControlBlock>) {
2      let future = TaskFuture::user_task(user_task.clone(), trap_loop(user_task));
3      executor::spawn(future);
4  }

```

代码 2-3 用户任务生成函数

## 2.4 多核心 CPU 管理

在 Del0n1x 中，我们将处理器抽象为 CPU，使用内核中 CPU 结构体进行统一管理。current 中保存当前正在运行的任务的 TCB；timer\_irq\_cnt 记录内核时钟中断次数，在内核时钟中断处理函数中会增加这个计数器，trap return 时会清零，

如果计数器大于阈值，手动对该任务 yield 进行调度，避免任务一直占用 CPU；kernel\_trap\_ret\_value 用于记录 pagafault 返回值。

```
1 pub struct CPU {  
2     current: Option<Arc<TaskControlBlock>>,  
3     timer_irq_cnt: usize,  
4     hart_id: usize,  
5     kernel_trap_ret_value: Option<SysResult<>>,  
6 }
```

代码 2-4 CPU 结构体定义

单个 CPU 被存放在全局的 PROCESSORS 管理器中，并对外暴露接口，通过管理器我们能获取到当前任务的上下文信息和页表 token、CPU id 号等。

```
1 const PROCESSOR: CPU = CPU::new();  
2 pub static PROCESSORS: SyncProcessors = SyncProcessors(UnsafeCell::new([PROCESSOR;  
    HART_NUM]));
```

代码 2-5 全局 CPU 管理器

## 2.5 任务控制块

进程是操作系统中资源管理的基本单位，而线程是操作系统中调度的基本单位。由于在 linux 设计理念中，线程是轻量级进程，所以在 Del0n1x 中使用统一的任务控制块来管理进程和线程。

```

1 pub struct TaskControlBlock {
2     pub pid: Pid, // 任务标识符
3     pub tgid: AtomicUsize, // leader 的 pid 号
4     pub pgid: AtomicUsize, // 进程组 id
5     pub task_status: SpinNoIrqLock<TaskStatus>, // 任务状态
6     pub thread_group: Shared<ThreadGroup>, // 线程组
7     pub memory_space: Shared<MemorySpace>, // 地址空间
8     pub parent: Shared<Option<Weak<TaskControlBlock>>>, // 父进程
9     pub children: Shared<BTreeMap<usize, Arc<TaskControlBlock>>>, // 子进程
10    pub fd_table: Shared<FdTable>, // 文件描述表
11    pub current_path: Shared<String>, // 路径
12    pub robust_list: Shared<RobustList>, // 存储线程的信息
13    pub futex_list: Shared<FutexBucket>, // futex 互斥锁队列
14    pub itimers: Shared<[ITimerVal; 3]>, // 任务的内部时钟
15    pub fsz_limit: Shared<Option<RLimit64>>, // 任务的资源限制
16    pub shmid_table: Shared<ShmidTable>,
17    pub pending: AtomicBool, // 是否有信号待处理
18    pub ucontext: AtomicUsize, // 信号用户态指针
19    pub sig_pending: SpinNoIrqLock<SigPending>, // 信号列表
20    pub blocked: SyncUnsafeCell<SigMask>, // 任务阻塞信号
21    pub handler: Shared<SigStruct>, // 信号处理集合
22    pub sig_stack: SyncUnsafeCell<Option<SignalStack>>, // 信号栈
23    pub waker: SyncUnsafeCell<Option<Waker>>, // 任务唤醒句柄
24    pub trap_cx: SyncUnsafeCell<TrapContext>, // 上下文
25    pub time_data: SyncUnsafeCell<TimeData>, // 时间戳
26    pub clear_child_tid: SyncUnsafeCell<Option<usize>>, // CHILD_CLEARID 清除地址
27    pub set_child_tid: SyncUnsafeCell<Option<usize>>, // CHILD_SETTID 设置地址
28    pub cpuset: SyncUnsafeCell<CpuSet>, // CPU 亲和性掩码
29    pub prio: SyncUnsafeCell<SchedParam>, // 调度优先级和策略
30    pub exit_code: AtomicI32, // 退出码
31 }

```

代码 2-6 任务控制块结构体

利用 rust Arc 引用计数和 clone 机制，可以有效的解决进程和线程之间资源共享和隔离问题。对于可以共享的资源，调用 `Arc::clone()` 仅增加引用计数（原子操作），未复制底层数据，父子进程共享同一份数据。如果是可以独立的资源（如 `memory_space`），调用 `clone` 会递归复制整个结构，生成完全独立的数据副本，父子进程修改互不影响。

### 2.5.1 进程和线程联系

在 Del0n1x 中，我们使用 `ThreadGroup` 管理线程组。其中选择 `BTreeMap` 作为管理数据结构，其中 key 是 task pid，value 是 TCB 的弱引用。线程之间需要一个 leader，而该 leader 是一个进程，同样被 `ThreadGroup` 管理。线程组的 leader 可以

用 `tgid` 表示；利用 `tgid` 可以通过 `BTreeMap` 快速定位到 `leader`。线程与线程之间并没有父子关系，他们同属于一个进程创建。

Type	Characteristic	Implementation
Process Leader	<code>tgid == pid</code>	新的内存空间，独立的信号处理程序
Thread	<code>tgid != pid</code>	通过 <code>CLONE_VM</code> 共享内存空间，通过 <code>CLONE_SIGHAND</code> 共享信号处理程序

进程和进程之间是树状结构，通过 `parent` 和 `children` 字段指明父进程和子进程。

Del0n1x 使用 `Manager` 管理任务和进程组，结构设计如下：

```

1  pub struct Manager {
2      pub task_manager: SpinNoIrqLock<TaskManager>,
3      pub process_group: SpinNoIrqLock<ProcessGroupManager>,
4  }
5
6  /// 存放所有任务的管理器，可以通过 pid 快速找到对应的 Task
7  pub struct TaskManager(pub HashMap<Pid, Weak<TaskControlBlock>>);
8  /// 存放进程组的管理器，通过进程组的 leader 的 pid 可以定位到进程组
9  pub struct ProcessGroupManager(HashMap<PGid, Vec<Pid>>);

```

代码 2-7 任务与进程组管理结构体

### 2.5.2 任务的状态

在 `rcore` 的基础上，我们为任务在运行过程中设计了 4 种状态：

- Ready: 任务已准备好执行，等待调度器分配 CPU 时间片；
- Running: 任务正在 CPU 上执行指令；
- Stopped: 任务被暂停执行，但未被终止，收到 `SIGSTOP` 信号
- Zombie: 任务已终止，但尚未被父进程回收

```

1  pub enum TaskStatus {
2      Ready,
3      Running,
4      Stopped,
5      Zombie,
6  }

```

代码 2-8 任务状态枚举

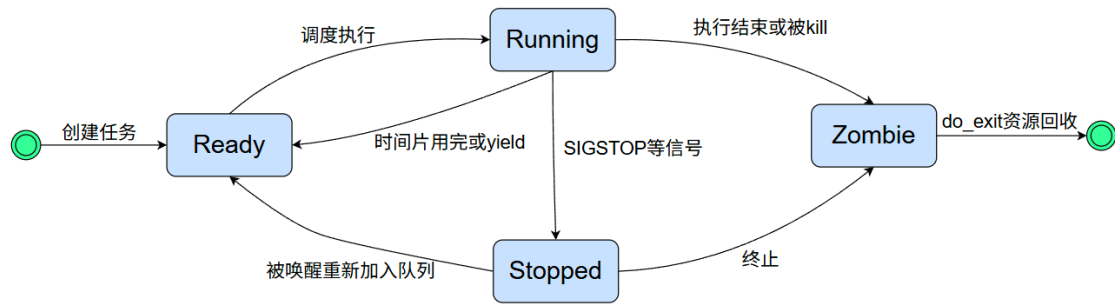


图 2-2 状态转化图

## 第 3 章 中断与异常处理

### 3.1 特权级切换

中断与异常是用户态与内核态之间切换的中重要机制。在 Del0n1x 中，为了更加清晰的设计模式，将中断异常分为两类：

- 用户中断异常：发生在用户态(U-Mode)的中断或异常
- 内核中断异常：发生在内核态(S-Mode)的中断或异常

当用户态发生中断或异常时，系统需要完成从用户态到内核态的切换。这个过程通过汇编函数 `__trap_from_user` 实现，这是用户态中断处理的入口点，负责保存完整的用户上下文。而 `sepc (era)`、`stval (estat)` 等寄存器则是由硬件自动完成保存。用户上下文保存在如下结构中：

```

1  pub struct TrapContext {
2      /* 0-31 */ pub user_gp: GPRs,
3      /* 32 */ pub sstatus: Sstatus,
4      /* 33 */ pub sepc: usize,
5      /* 34 */ pub kernel_sp: usize,
6      /* 35 */ pub kernel_ra: usize,
7      /* 36-47 */ pub kernel_s: [usize; 12],
8      /* 48 */ pub kernel_fp: usize,
9      /* 49 */ pub kernel_tp: usize,
10     /* 50 */ pub float_regs: UserFloatRegs,
11 }

```

代码 3-1 TrapContext 结构

需要注意的是，在 riscv64 中，我们通过检测 `fs` 寄存器是否被使用决定是否保存浮点寄存器，这样的设计减少上下文切换的开销。但是通过查阅手册，loongarch64 架

构中并没有提供这样的寄存器，所以我们将 loongarch64 中浮点寄存器的保存处理放置在 `__trap_from_user`，同时我们在 `TrapContext` 中增加了 `fcsr` 字段用于保存浮点控制状态寄存器。`fcsr` 是浮点运算单元（FPU）的核心控制寄存器，它负责管理浮点运算的异常标志、舍入模式、使能控制等关键功能。

## 3.2 处理过程

### 3.2.1 内核中断异常处理

以下均以 riscv64 为例。Del0n1x 的内核态中断异常处理目前支持时钟中断和地址缺页异常（对于外部中断放置于比赛下一阶段进行完善），代码如下：

```

1  Trap::Interrupt(Interrupt::SupervisorTimer) => {
2      TIMER_QUEUE.handle_expired();
3      get_current_cpu().timer_irq_inc();
4      set_next_trigger();
5  }
6  Trap::Exception(e) => match e {
7      Exception::StorePageFault
8      | Exception::InstructionPageFault
9      | Exception::LoadPageFault => {
10     let access_type = match e {
11         Exception::InstructionPageFault => PageFaultAccessType::RX,
12         Exception::LoadPageFault => PageFaultAccessType::RO,
13         Exception::StorePageFault => PageFaultAccessType::RW,
14         _ => unreachable!(),
15     };
16
17     let task = current_task().unwrap();
18     result = task
19         .with_mut_memory_space(|m| {
20             m.handle_page_fault(stval.into(), access_type)
21         });
22     result.is_err().then(|| {
23         use crate::hal::arch::current_inst_len;
24         sepc::write(sepc + current_inst_len());
25     });
26 }
27
28 set_ktrap_ret(result);

```

代码 3-2 内核中断异常处理关键函数

内核缺页异常通常发生在系统调用中，由于进程虚拟地址空间没有实际分配独立的物理页帧，而是将其指向父进程对应的物理页帧，当内核向用户传入的地址写入数据时，会触发 `PageFault` 异常跳转至内核处理函数中，在这次 `trap` 的处理中，我们

不仅对子进程分配实际的物理页帧，恢复相应页表项的标志位，同时通过检查 trap 的返回值实现用户地址空间可写性的检查。

### 3.2.2 用户中断异常处理

在 Del0n1x 中除了对用户态的 PageFault 处理之外，还实现了时钟中断和系统调用处理。对于时钟中断，Del0n1x 检查了全局定时器中是否有超时任务，然后设置下一次时钟中断时间点，最后需要调用 yield 释放当前任务对 CPU 的使用权，调度下一个任务，避免任务长时间占用 CPU 导致其他任务饥饿。对于系统调用处理，会先将中断上下文中的 sepc 加 4，使得从内核态返回到用户态后能够跳转到下一条指令。然后，调用 syscall 函数系统调用。系统调用完成后，将返回值保存在 x10 寄存器。

```

1  ....
2  Trap::Exception(Exception::UserEnvCall) => { // 7
3      let mut cx = current_trap_cx();
4      let old_sepc: usize = cx.get_sepc();
5      let syscall_id = cx.user_gp.a7;
6      cx.set_sepc(old_sepc + 4);
7
8      let result = syscall(
9          syscall_id,
10         [cx.user_gp.a0,
11          cx.user_gp.a1,
12          cx.user_gp.a2,
13          cx.user_gp.a3,
14          cx.user_gp.a4,
15          cx.user_gp.a5]
16     ).await;
17
18     // cx is changed during sys_exec, so we have to call it again
19     cx = current_trap_cx();
20
21     match result {
22         Ok(ret) => {
23             cx.user_gp.a0 = ret as usize;
24         }
25         Err(err) => {
26             if (err as isize) < 0 {
27                 cx.user_gp.a0 = err as usize;
28             } else {
29                 cx.user_gp.a0 = (-(err as isize)) as usize;
30                 info!("[syscall ret] sysID = {}, errmsg: {}", syscall_id, err.get_info());
31             }
32         }
33     }
34 }
35
36 Trap::Interrupt(Interrupt::SupervisorTimer) => { // 5
37     TIMER_QUEUE.handle_expired();
38     set_next_trigger();
39     yield_now().await;
40 }
41 ....

```

代码 3-3 用户中断异常处理函数

### 3.2.3 返回用户态

Del0n1x 中从内核态返回到用户态过程交付予 `user_trap_return` 函数处理。在该函数的处理逻辑中，首先要通过设置 `stvec` 寄存器确保下次用户中断的入口地址正确，然后恢复用户浮点寄存器状态，同时修改进程时间戳记录进程 `trap_out` 时间。最后通过 `__return_to_user` 调用 `sret` 实现 S 监督模式到 U 用户模式特权级的切换。



```

1 pub fn user_trap_return() {
2     // 重新修改 stvec 设置 user 的 trap handler entry
3     set_trap_handler(IndertifyMode::User);
4
5     let trap_cx = current_trap_cx();
6     trap_cx.float_regs.trap_out_do_with_freg();
7     trap_cx.sstatus.set_fs(FS::Clean);
8
9     get_current_cpu().timer_irq_reset();
10    let task = current_task().unwrap();
11    task.get_time_data_mut().set_trap_out_time();
12    unsafe {
13        __return_to_user(trap_cx);
14    }
15    task.get_time_data_mut().set_trap_in_time();
16
17    trap_cx.float_regs.trap_in_do_with_freg(trap_cx.sstatus);
18 }

```

代码 3-4 内核-&gt;用户切换

## 第 4 章 内存管理

### 4.1 物理内存管理

#### 4.1.1 内核动态内存分配器

Phoenix 使用伙伴分配器管理内核所需的动态内存结构，来自 `crate buddy_system_allocator`。

伙伴分配器 (Buddy Allocator) 是一种内存分配算法，常用于操作系统内核和高性能应用程序中，通过分配和管理内存块来满足不同大小的内存请求，并进行高效的合并和分割操作。其工作原理是将内存块按照 2 的幂次方大小分为多个层级，当需要分配特定大小的内存时，从最小适合该请求的层级开始查找。每个内存块都有一个“伙伴”块，如果块大小是  $2^k$ ，那么它的伙伴块也是  $2^k$  且紧挨着它。通过检查和计算伙伴块的地址，可以快速地进行内存分割与合并。当需要分配的内存块小于当前可用最小块时，将当前块一分为二，直到找到合适大小的块为止；当释放一个内存块时，若其伙伴块也空闲，则将两个块合并为一个更大的块，递归进行直到不能再合并为止。伙伴分配器的优点在于分配和释放内存块的操作非常快速，且通过内存块大小的选择和合并操作，有效减少了外部碎片。

### 4.1.2 物理页分配器

内核还需要管理全部的空闲物理内存，Phoenix 为此使用了来自 rCore 的仓库的 bitmap-allocator。Phoenix 在内核初始化时，会将所有内核未占用的物理内存加入物理页分配器。

Bitmap allocator 的主要原理是通过一个位图来管理一段连续的内存空间。这个位图中的每一位代表一块内存，如果该位为 0，说明对应的内存块空闲；如果该位为 1，说明对应的内存块已经被分配出去。当需要分配一个指定大小的内存时，bitmap allocator 首先检查位图中是否有足够的连续空闲内存块可以满足分配请求。如果有，就将对应的位图标记为已分配，并返回该内存块的起始地址；如果没有，就返回空指针，表示分配失败。当需要释放已经分配出去的内存时，bitmap allocator 将对应位图标记为未分配。这样，已经释放的内存块就可以被下一次分配请求使用了。

此外，Phoenix 将物理页帧抽象成 FrameTracker 结构体，并结合 RAII 的思想，在结构体析构时自动调用 dealloc\_frame 函数将页帧释放。

```

1  /// Manage a frame which has the same lifecycle as the tracker.
2  pub struct FrameTracker {
3      /// PPN of the frame.
4      pub ppn: PhysPageNum,
5  }
6
7  impl Drop for FrameTracker {
8      fn drop(&mut self) {
9          dealloc_frame(self.ppn);
10     }
11 }
```

代码 4-1 FrameTracker 结构体

## 4.2 地址空间

### 4.2.1 地址空间布局

Phoenix 地址空间的设计如下图所示：

Phoenix 内核态页表保存在全局内核地址空间 KERNEL\_SPACE 中，用户地址空间共享内核二级页表。

对于内核地址空间,为了方便管理所有物理地址,采用偏移映射的方式将物理地址以加上偏移量的方式映射为虚拟地址,即每一个虚拟地址都为对应物理地址加上 `VIRT_RAM_OFFSET`。

对于用户地址空间,为了利用分页机制的灵活性,消除外部碎片,采用随机映射的方式将需要的页随机映射到空闲物理内存中随机一块页帧。

#### 4.2.2 Boot 阶段高位映射

在 QEMU 平台上,内核的入口地址位于 `0x8020_0000`。在 `xv6` 与 `rCore-Tutorial` 中,`0x8020_0000` 以上部分以直接映射的方式作为内核地址空间,`0x8000_0000` 以下部分以随机映射的方式作为用户程序地址空间。因此用户程序最多只有 2G 的地址空间,而考虑到 MMIO 也映射在低位,用户程序地址空间只会更少。因此,Phoenix 充分利用页表的灵活性,将内核地址空间以偏移映射的方式映射到 RISC-V SV39 规定的高位地址空间,即 `0xFFFF_FFC0_0000_0000` 至 `0xFFFF_FFFF_FFFF_FFFF`,每一个虚拟地址对应于物理地址加上一个相同的偏移量,而将用户地址空间映射到 `0x0` 至 `0x3F_FFFF_FFFF`。这样,用户地址空间不足的问题就被解决了。

但是问题也随之而来,OpenSBI 会识别内核 ELF 文件入口地址的加载地址 (LMA),然后在启动完毕后会跳转到此处,这时 RISC-V 页表机制还未开启,内核会直接访问物理地址。而开启页表时需要保证指令在启动页表前后物理地址连续,这就需要跳板页。Phoenix 希望在内核尽早映射到高位地址空间,这样在调试时能够尽早将代码与地址对应,因此 Phoenix 首先将内核虚拟地址链接到高位地址空间,将加载地址链接到低位,然后结合 RISC-V 巨页的机制,在内核启动阶段,在 Boot 页表中构造了三个巨页,地址分别为 `0x8000_0000`、`0xFFFF_FFC0_0000_0000` 和 `0xFFFF_FFC0_8000_0000`。其中,`0x8000_0000` 所在巨页作为跳板页,保证在打开页表前后 `pc` 寄存器指向连续的物理地址;`0xFFFF_FFC0_0000_0000` 所在巨页保存了 MMIO 的高位映射,以便在 Boot 页表阶段开启打印调试功能;`0xFFFF_FFC0_8000_0000` 所在巨页对应于内核 ELF 文件的虚拟地址高位映射,在 `_start` 函数最后会跳转到此处执行内核代码。在 `entry.rs` 代码执行完毕后,内核便成功执行在高位地址空间了,最后,Boot 页表会在 `main` 函数执行过程中被更加细化的全局内核页表 `KERNEL_SPACE` 所取代。

```

1  #[link_section = ".bss.stack"]
2  static mut BOOT_STACK: [u8; KERNEL_STACK_SIZE * MAX_HARTS] =
3      [0u8; KERNEL_STACK_SIZE * MAX_HARTS];
4
5  #[repr(C, align(4096))]
6  struct BootPageTable([u64; PTES_PER_PAGE]);
7
8  static mut BOOT_PAGE_TABLE: BootPageTable = {
9      let mut arr: [u64; PTES_PER_PAGE] = [0; PTES_PER_PAGE];
10     arr[2] = (0x80000 << 10) | 0xcf;
11     arr[256] = (0x00000 << 10) | 0xcf;
12     arr[258] = (0x80000 << 10) | 0xcf;
13     BootPageTable(arr)
14 };
15
16 unsafe extern "C" fn _start(hart_id: usize, dtb_addr: usize) -> ! {
17     core::arch::asm!(
18         // 1. set boot stack
19         "
20             addi    t0, a0, 1
21             slli    t0, t0, 16          // t0 = (hart_id + 1) * 64KB
22             la      sp, {boot_stack}
23             add     sp, sp, t0          // set boot stack
24         ",
25         // 2. enable sv39 page table
26         "
27             la      t0, {page_table}
28             srli    t0, t0, 12
29             li      t1, 8 << 60
30             or      t0, t0, t1
31             csrw    satp, t0
32             sfence.vma
33         ",
34         // 3. jump to rust_main
35         "
36             li      t2, {virt_ram_offset}
37             or      sp, sp, t2
38             la      a2, rust_main
39             or      a2, a2, t2
40             jalr    a2                // call rust_main
41         ",
42         boot_stack = sym BOOT_STACK,
43         page_table = sym BOOT_PAGE_TABLE,
44         virt_ram_offset = const VIRT_RAM_OFFSET,
45         options(noreturn),
46     )
47 }

```

代码 4-2 Boot 阶段高位映射代码

#### 4.2.3 地址空间管理

Phoenix 使用 RAII 机制管理地址空间，主要使用 `MemorySpace` 和 `VmArea` 以及 `RangeMap` 数据结构。

```

1  /// Virtual memory space for kernel and user.
2  pub struct MemorySpace {
3      /// Page table of this memory space.
4      page_table: PageTable,
5      /// Map of `VmArea`s in this memory space.
6      areas: RangeMap<VirtAddr, VmArea>,
7  }
8
9  /// A contiguous virtual memory area.
10 pub struct VmArea {
11     /// Aligned `VirtAddr` range for the `VmArea`.
12     range_va: Range<VirtAddr>,
13     /// Hold pages with RAII.
14     pub pages: BTreeMap<VirtPageNum, Arc<Page>>,
15     /// Map permission of this area.
16     pub map_perm: MapPerm,
17     /// Type of this area.
18     pub vma_type: VmAreaType,
19
20     // For mmap.
21     /// Mmap flags.
22     pub mmap_flags: MmapFlags,
23     /// The underlying file being mapped.
24     pub backed_file: Option<Arc<dyn File>>,
25     /// Start offset in the file.
26     pub offset: usize,
27 }
28
29 /// A range map that stores range as key.
30 pub struct RangeMap<U: Ord + Copy + Add<usize>, V>(BTreeMap<U, Node<U, V>>);

```

代码 4-3 地址空间管理结构体

## 4.3 缺页异常处理

Phoenix 目前能够利用缺页异常处理来实现写时复制 (Copy on write) 懒分配 (Lazy page allocation) 以及用户地址检查机制和零拷贝技术。

当用户程序因缺页异常返回内核时，内核异常处理函数能够从 `stval` 寄存器读取异常发生的地址，并交给 `VmArea::handle_page_fault` 函数进行处理。

### 4.3.1 CoW 写时复制技术

在 `fork` 进程时，Phoenix 会将原 `MemorySpace` 中的除共享内存外每一个已分配页的 PTE 都删除写标志位，打上 COW 标志位，然后重新映射到页表中，并将。在用户向 COW 页写入时会触发缺页异常陷入内核，在 `VmArea::handle_page_fault` 函数中，

内核会根据 COW 标志位转发给 COW 缺页异常处理函数，缺页异常处理函数会根据 Arc<Page> 的原子持有计数判断是否为最后一个持有者，如果不是最后一个持有者，会新分配一个页并复制原始页的数据并恢复写标志位重新映射，如果是最后一个持有者，直接恢复写标志位。

```

1  impl MemorySpace {
2      /// Clone a same `MemorySpace` lazily.
3      pub fn from_user_lazily(user_space: &mut Self) -> Self {
4          let mut memory_space = Self::new_user();
5          for (range, area) in user_space.areas().iter() {
6              let mut new_area = area.clone();
7              for vpn in area.range_vpn() {
8                  if let Some(page) = area.pages.get(&vpn) {
9                      let pte = user_space
10                         .page_table_mut()
11                         .find_leaf_pte(vpn)
12                         .unwrap();
13                      let (pte_flags, ppn) = match area.vma_type {
14                          VmAreaType::Shm => {
15                              // no cow for shared memory
16                              new_area.pages.insert(vpn, page.clone());
17                              (pte.flags(), page.ppn())
18                          }
19                          _ => {
20                              // copy on write
21                              let mut new_flags = pte.flags() | PTEFlags::COW;
22                              new_flags.remove(PTEFlags::W);
23                              pte.set_flags(new_flags);
24                              (new_flags, page.ppn())
25                          }
26                      };
27                      memory_space.page_table_mut().map(vpn, ppn, pte_flags);
28                  } else {
29                      // do nothing for lazy allocated area
30                  }
31              }
32              memory_space.push_vma_lazily(new_area);
33          }
34          memory_space
35      }
36  }

```

代码 4-4 写时复制实现

#### 4.3.2 懒分配技术

懒分配技术主要用于堆栈分配以及 mmap 匿名映射或文件映射。在传统的内存分配方法中，操作系统在进程请求内存时会立即为其分配实际的物理内存。然而，

这种方法在某些情况下可能导致资源的浪费，因为进程可能并不会立即使用全部分配的内存。

懒分配技术的核心思想是推迟实际物理内存的分配，直到进程真正访问到该内存区域。这样可以优化内存使用，提高系统性能。

对于内存的懒分配，比如堆栈分配，mmap 匿名内存分配，Phoenix 将许可分配的范围记录下来，但并不进行实际分配操作，当用户访问到许诺分配但未分配的页面时会触发缺页异常，缺页异常处理函数会进行实际的分配操作。

对于 mmap 文件的懒分配，Phoenix 将其与页缓存机制深度融合，Phoenix 同样执行懒分配操作，当缺页异常时再从页缓存中获取页面。

#### 4.3.3 用户地址检查与零拷贝技术

在系统调用过程中，内核需要频繁与用户态指针指向的数据进行交互。在 Phoenix 中，用户和内核共享地址空间，因此在访问用户态的内存时不需要同 xv6 那样通过软件查询页表，而是可以利用硬件页表机制直接解引用用户态指针。

然而，用户态指针并不总是有效的，有可能指向非法内存，出于安全性保证，内核需要能够捕获这种异常。在用户态下，无效指针解引用，或者向只读地址写入数据会触发缺页异常，陷入内核并执行相应缺页异常处理函数，通常，缺页异常处理函数会向程序发送 SIGSEGV 信号或者终止进程。而 Phoenix 内核态也需要解引用用户态的指针，因此内核也需要捕获并处理这种页错误。

我们参考了往届 MankorOS 队伍的做法，借助硬件 MMU 部件的帮助实现了高效的指针检查。该做法的基本思路是，先将内核的异常捕捉函数替换为“用户检查模式”下的函数，然后直接尝试向目标地址读取或写入一个字节。若是目标地址发生了缺页异常，则内核将表现得如同用户程序发生了一次异常一般，进入用户缺页异常处理程序进行处理。若处理成功或目标地址访问成功，便可假定当前整个页范围内都是合法的用户地址空间，否则用户指针便不合法。该处理方法相当于直接利用了硬件 MMU 来检查用户指针是否可读或可写，在用户指针正常时速度极快，同时还能完全复用用户缺页异常处理的代码来处理用户指针懒加载/CoW 的情况。

此外，Phoenix 基于此实现了用户态指针内容的零拷贝技术，即内核态不需要软件模拟地址翻译并复制用户态指针指向的数据到内核态，而是可以直接访问用户态指针，避免用户态数据到内核态数据的拷贝。用户态传入的指针经过检查后，会被转换成 UserRef、UserMut 和 UserSlice 对象。每一个对象都存放具体的指针，并保存了 SumGuard 以获取内核态访问用户地址空间的权限。Phoenix 充分利用了

Rust 提供的类型机制，为上述对象实现了 `deref` 方法，使得在不改变外部函数签名的情况下，依然能保持 `sstatus` 寄存器 `sum` 位开启，直接访问用户地址空间。这种实现不仅高效，还极大缩短了代码量。例如，`sys_read` 函数只需短短 3 行代码，不仅实现了用户地址空间检查，还能将文件内容零拷贝直接填充到用户提供的缓冲区。

```
/// User slice. Hold slice from `UserPtr` and a `SumGuard` to provide user
/// space access.
pub struct UserSlice<'a, T> {
    slice: &'a mut [T],
    _guard: SumGuard,
}

impl<'a, T> core::ops::DerefMut for UserSlice<'a, T> {
    fn deref_mut(&mut self) -> &mut Self::Target {
        &mut self.slice
    }
}

pub async fn sys_read(
    &self,
    fd: usize,
    buf: UserWritePtr<u8>,
    count: usize,
) -> SyscallResult {
    let file = self.task.with_fd_table(|table| table.get_file(fd))?;
    let mut buf: UserSlice<'_, u8> = buf.into_mut_slice(&task, count)?;
    file.read(&mut buf).await
}
```

## 第 5 章 文件系统

### 5.1 虚拟文件系统

虚拟文件系统（Virtual File System，简称 VFS）是对各种文件系统的抽象，这种抽象屏蔽了各种具体文件系统的细节，为内核提供统一的统一的文件系统接口。



在标准的系统调用中，例如 `open()`、`read()`、`write()` 中，可以忽略掉文件系统的特性，调用统一的 VFS 接口，从而简洁、安全地实现对文件的各种操作。

我们充分利用 rust 语言的特性，借鉴 Linux 文件系统设计，以面向 trait (中文翻译为属性) 的方式进行编程，提供了方便、安全的文件系统接口供内核的其他模块与用户系统调用使用。

Del0n1x OS 的虚拟文件系统的主要数据结构为 SuperBlock、Inode、Dentry、File。

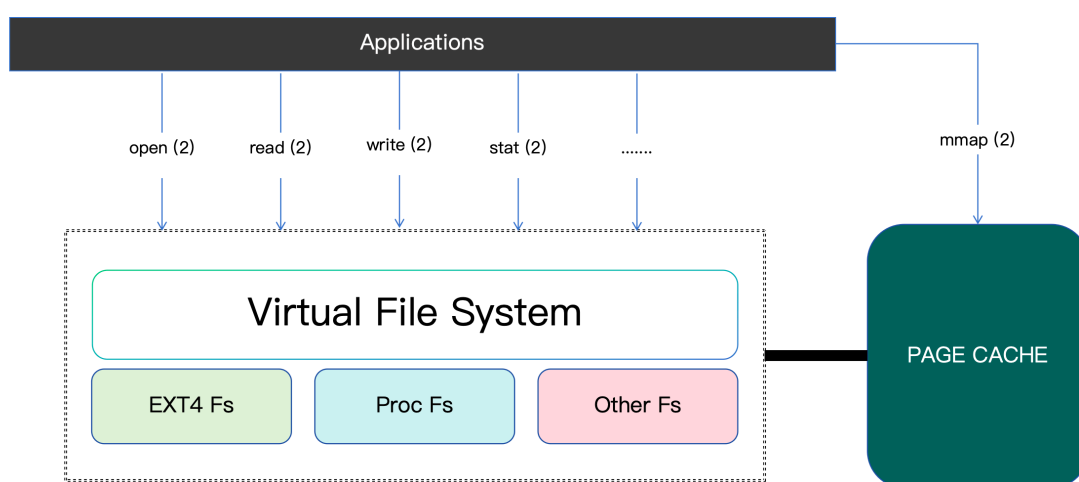


图 5-1 虚拟文件系统

### 5.1.1 SuperBlock

SuperBlock trait 超级块属性是一个具体的文件系统的抽象。每一个具体文件系统都有对 SuperBlock trait 的实现，不同文件系统对同一个方法会有不同的实现。

超级块 SuperBlock trait 的定义如下：

```
1 pub trait SuperBlockTrait: Send + Sync {
2     /// 获取根节点
3     fn root_inode(&self) -> Arc<dyn InodeTrait>;
4     /// 将数据写回
5     fn sync(&self);
6     /// 显示文件系统信息
7     fn fs_stat(&self) -> StatFs;
8 }
```

代码 5-1 文件系统超级块

与 C 语言相比 ,rust 提供了足够的抽象机制。对于不同的文件系统 ,实现 SuperBlock trait 后就可以在文件系统中安全和高效地使用 ,鲜明的方法名称也为编程带来方便。rust 同时区别于传统的面向对象语言 ,抛弃了继承机制的设计 ,鼓励面向 trait (翻译为属性) 编程 ,并使用组合替代继承 ,使得代码结构更为简单高效。

我们的内核中简化了超级块的提供的功能 ,仅提供有限的信息

#### 5.1.2 Inode

索引节点 (inode) 是文件系统的核心 ,是文件的抽象。

索引节点属性 InodeTrait 是对索引节点 (inode) 的抽象 ,不同的文件系统有其自身的实现 ,但是其暴露出的方法是统一的。Linux 当中 inode 结构体有 struct inode\_operations \*i\_op 字段 ,InodeTrait 与索引节点 (inode) 的 inode\_operations 功能类似 ,定义了索引节点 (inode) 应当实现的功能。

不同文件系统的 inode 结构体应当实现 InodeTrait。在 Linux 中 ,当处理内核文件系统操作或文件系统相关系统调用的时候 ,会获得对应的 inode 对象 ,随后去检查并调用 inode->i\_op 中的回调函数来实现具体功能。这种面向方法 (也可以成为面向对象) 的编程思想极大地提高了内核编码的安全性和便利性 ,我们无需关心各个文件系统当中如何实现具体的方法的 ,我们仅仅需调用暴露的方法就好了。

rust 可以很方便地做到这一点 ,因为 rust 原生面向 Trait 编程。我们充分利用了这一点 !

我们定义了索引节点 (Inode) 应当实现的方法 InodeTrait ,具体的文件系统需要实现其中的方法 ,或者使用默认方法 ,就可以供上层使用

索引节点属性 InodeTrait 的主要内容如下 :

```

1 pub trait InodeTrait: Send + Sync {
2     /// inode 的信息
3     fn fstat(&self) -> Kstat;
4     /// 在文件夹上创建一个子文件
5     fn do_create(&self, bare_dentry: Arc<Dentry>, _ty: InodeType) -> Option<Arc<dyn
6     InodeTrait>>;
7     /// 读文件
8     fn read_at(&self, _off: usize, _buf: &mut [u8]) -> usize;
9     /// 写文件
10    fn write_at(&self, _off: usize, _buf: &[u8]) -> usize;
11    /// 截断文件
12    fn truncate(&self, _size: usize) -> usize;
13    /// unlink 文件
14    fn unlink(&self, valid_dentry: Arc<Dentry>) -> SysResult<usize>;
15    /// link 文件
16    fn link(&self, bare_dentry: Arc<Dentry>) -> SysResult<usize>;
17    /// 获得 page cache
18    fn get_page_cache(&self) -> Option<Arc<PageCache>>;
19    /// 更名
20    fn rename(&self, old_path: Arc<Dentry>, new_path: Arc<Dentry>) -> SysResult<usize>;
21    /// 文件夹读取目录项
22    fn read_dents(&self) -> Option<Vec<Dirent>>;
23    /// io 操作
24    fn ioctl(&self, op: usize, arg: usize) -> SysResult<usize>
25 }

```

代码 5-2 InodeTrait 接口定义

在文件对象 File 或者目录项对象中会持有数据类型为 `Arc<dyn InodeTrait>` 的索引节点 (inode) 对象。通过索引节点 (inode) 对象, 通过其实现的 Inode trait 接口我们可以获得文件的各种信息, 进行文件的读写、创建、删除等操作。在调用 InodeTrait 定义的方法的时候, 我们不需要关心其具体的数据类型, 这一行为会自动分派给对应的文件系统的索引节点 (inode) 实现。

### 5.1.3 Dentry

目录项 Dentry 是目录树上节点的抽象。每一个有效的目录项持有对一个合法的索引节点 (inode) 的引用操作系统应用程序使用路径获得对应的文件, 这一过程是由目录项 (dentry) 完成。目录项构成了一颗目录树, 一般而言, 通过从挂载节点向下搜索, 实现对文件的查找。同时 Dentry 还对文件系统目录树进行了管理。

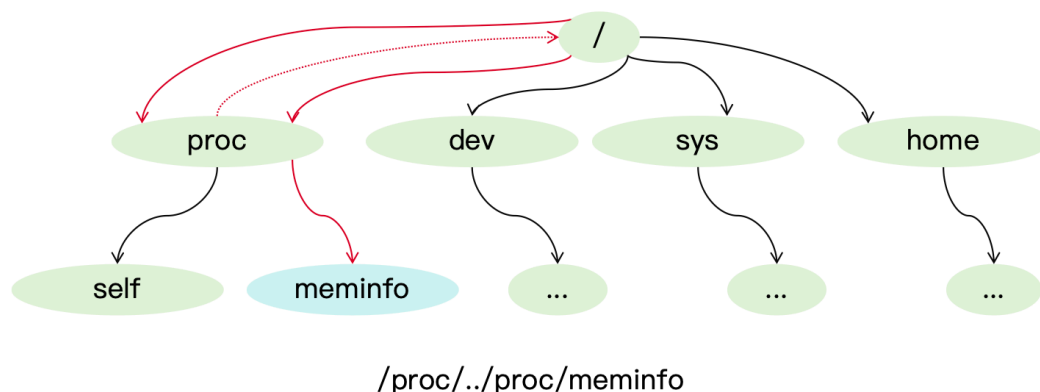


图 5-2 访问目录树

在初始化目录树的过程中，我们需要将具体的文件系统挂载（mount）在目录树上，目录树提供了这一功能的实现。我们可以把任意的文件系统挂载到目录树上的文件夹节点上，将该目录项（dentry）持有的 inode 替换为该文件系统的根目录，我们就可以通过对应的路径，去访问该文件系统下的文件。

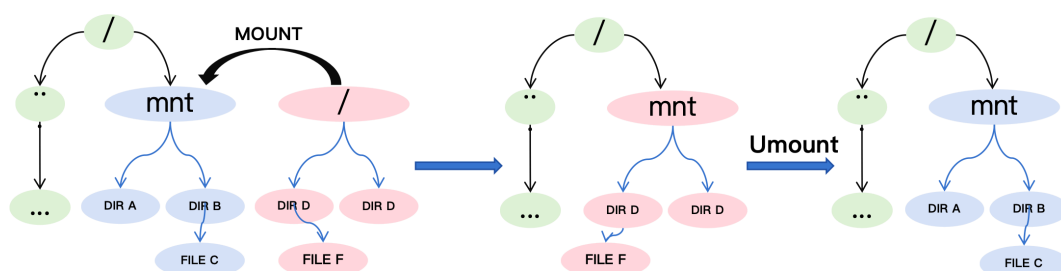


图 5-3 mount and umount

在我们的实现中，当我们挂载了一个新的文件系统到某个目录项上时，原有的目录项上持有的索引节点（inode）的引用会隐藏起来，当从这个目录项上卸载（unmount）该文件系统的时候原来持有的索引节点（inode）会恢复，进而恢复原有的子树结构。

目前仅支持内核中调用，未来我们希望实现 loop 设备后在用户态也可以使用目录项 Dentry 的定义如下：

```
pub struct Dentry {
    /// 目录项文件名
    name: RwLock<String>,
```

```

    /// 对父 dentry 的弱引用
    parent: Weak<Dentry>,
    /// 孩子 dentry 的强引用
    children: RwLock<HashMap<String, Arc<Dentry>>>>,
    /// 当前的持有的 inode 对象
    inode: RwLock<Vec<Arc<dyn InodeTrait>>>>,
    /// dentry 的状态
    status: RwLock<DentryStatus>,
}

```

通过 children 字段获得当前目录项的子目录项,通过 parent 获得当前目录项的双亲目录项,注意到这里使用弱引用(不增加引用计数)防止出现循环引用。inode 字段为所持有的索引节点(inode)。

定义目录项状态 DentryStatus

```

1  pub enum DentryStatus {
2      /// 这个 dentry 是有效的,并且已经初始化
3      Valid,
4      /// 这个 dentry 是有效的,但是没有初始化
5      Unint,
6      /// 这个 dentry 是无效的
7      Negtive,
8  }

```

代码 5-3 DentryStatus 枚举

目录项状态在这三个状态之间转移,当目录项被标记为无效时,会在合适的时机进行回收,并且释放对索引节点(inode)对象的引用。当目录项尚未初始化的时候,会在访问时初始化。只有当目录项有效时才可以进行访问。

目录项(dentry)额外使用了一个缓存(Cache)用于加速从路径到目录项的查找,目录项缓存(DentryCache)使用内核定义的Cache泛型容器进行定义。目录项缓存的存在极大地加速了获得索引节点(inode)的过程,获得显著的性能提升。

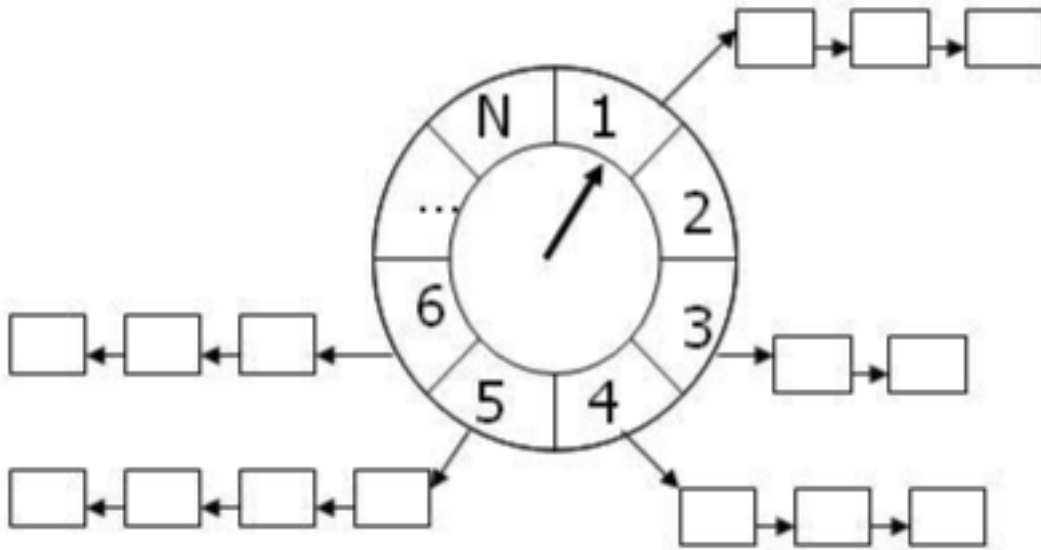


图 5-4 没有 cache，有 cache 在 musl 和 glib 下的性能对比

#### 5.1.4 File

文件对象 (file) 是进程中已打开文件在内核中的表示。文件由系统调用 `open()` 创建，由系统调用 `close()` 关闭。一个最基本的文件对象 (file) 由持有的索引节点 (inode)、文件的偏移和其他状态信息组成。用户态程序调用 `open` 系统调用来创建文件对象时，首先需要根据路径获得对应的目录项 (dentry)，通过目录项获得索引节点 (inode)，将索引节点包装在文件对象 (file) 中，并设置其状态，最后将创建好的文件对象注册到文件描述符表。

## 5.2 磁盘文件系统

磁盘文件系统 (Disk File System) 是操作系统中用于管理磁盘和其他存储设备上数据存储和组织的机制。其定义了数据如何以文件和目录的形式存储、命名、访问、更新，以及如何在物理存储介质上分配空间。

### 5.2.1 EXT4 文件系统

Ext4 (第四代扩展文件系统) 是 Ext3 文件系统的继承者，主要用于 Linux 操作系统。与前代文件系统相比，Ext4 在性能、可靠性和容量方面都有显著改进。Ext4 文件系统对 Unix 操作系统适配性更好，支持硬链接等操作。

我们的操作系统使用了开源的 `lwext4-rust` 库，为其提供可用的持久化设备，并利用库提供的功能，根据 VFS 的接口设计进行适配与抽象。

## 5.3 非磁盘文件系统

非磁盘文件系统 ( Non-Disk File System ) 是指不依赖于传统磁盘存储介质的文件系统。在我们的操作系统中提供了若干个非磁盘文件系统供系统使用。

### 5.3.1 procfs

procfs 是一种特殊的文件系统，它不是从磁盘上的文件系统中读取数据，而是从内核中读取数据。procfs 包括：

- /proc/mounts : 显示当前挂载的文件系统
- /proc/meminfo : 提供关于系统内存使用情况的信息，包括总内存、可用内存、缓存和缓冲区等详细数据
- /proc/exe : 当前正在运行的程序
- /proc/self : 当前正在运行的进程所持有的内容

这个文件系统完整地实现了 VFS 中所有的接口，用户可以透明地使用其中的文件。用户态程序可以很方便地从这些文件中提取相关信息。

### 5.3.2 devfs

devfs 中的文件代表一些具体的设备，比如终端、硬盘等。devfs 内包含：

- /dev/zero : 一个无限长的全 0 文件
- /dev/null : 用于丢弃所有写入的数据，并且读取时会立即返回 EOF ( 文件结束 )
- /dev/random : 一个伪随机数生成器，提供随机数据流
- /dev/rtc : 实时时钟设备，提供日期和时间
- /dev/tty : 终端设备，能支持 ioctl 中的特定命令
- /dev/loop0 : 回环设备，用于虚拟块设备

## 5.4 页缓存

页缓存 ( Page Cache ) 以页为单位缓存文件的内容。当我们需要读文件的时候，在缓存命中的情况下，就不需要去访问持久化设备，从而提高性能。当写文件的时候，也可以直接往页缓存中写，同时标记为脏页，就可以直接返回，而不需要等待数据真正地被写入到磁盘。脏页由内核进行统一的管理。总而言之，页缓存的设计极大地提高了文件的读写性能。

页缓存同时也是连接文件系统模块和内存模块的桥梁。用户可以调用 `mmap` 系统调用，将文件的页缓存映射到用户态地址空间中，用户就可以安全且高效地实现对文件的访问。`mmap` 设计可以复用页缓存机制而安全高效地实现

在用户态程序借助文件进行进程间通信（IPC）这一常见场景下，性能也会获得极大地提升。

以下为页缓存（Page Cache）的定义

```
1 pub struct PageCache {
2     pub pages: RwLock<BTreeMap<usize, Arc<Page>>>,
3     inode: RwLock<Option<Weak<dyn InodeTrait>>>,
4 }
```

代码 5-4 PageCache 结构体

在 PageCache 实现中以页对齐的地址为 key 去获得对应的页帧（Page）。

Ext4 文件系统的索引节点（inode）会持有这一页缓存，当不做特别要求的时候，均在页缓存中读写。系统会对页缓存进行管理，实现对齐其创建，映射与释放。并且当索引节点（inode）被释放的时候，其持有的页缓存（PageCache）也会自动释放，进而释放所持有的页帧。

## 5.5 其他数据结构

### 5.5.1 FdTable

用户态程序使用 `open` 系统调用打开文件后，会获得文件描述符（File Descriptor）来用于控制文件（File），这需要内核为每一个进程创建一个对应的文件描述符表（Fd Table）用于实现文件描述符转换到文件的映射

以下为文件描述符表的实现：

```
pub struct FdTable {
    pub table: Vec<FdInfo>,
    pub rlimit: RLimit64,
    free_bitmap: Vec<u64>,
    next_free: usize,
    freed_stack: Vec<usize>,
}

#[derive(Clone)]
```



```
pub struct FdInfo {  
    pub file: Option<Arc<dyn FileTrait>>,  
    pub flags: OpenFlags,  
}
```

## 第 6 章 进程间通信

### 6.1 信号机制

信号是操作系统向进程传递事件通知的一种机制，主要用于通知进程发生了异步事件。在我们的内核中，严格按照 Linux 中对于信号结构的设计，实现了相对完善且清晰的信号机制。信号相关结构体设计自顶向下分别为，`SigStruct`（一个包含所有信号处理方法的数组），其中每个元素为 `KSigAction`（内核层信号动作），`SigAction`（信号处理相关配置），三者关系如下：

```

1  #[derive(Clone, Copy)]
2  pub struct SigStruct {
3      pub actions: [KSigAction; MAX_SIGNUM],
4  }
5
6  /// 内核层信号动作
7  #[derive(Clone, Copy)]
8  pub struct KSigAction {
9      pub sa: SigAction,
10     pub sa_type: SigHandlerType,
11 }
12
13 /// 用户层信号处理配
14 #[derive(Clone, Copy, Debug)]
15 #[repr(C)]
16 pub struct SigAction {
17     /// 信号处理函数类型，可能是自定义，也可能是默认
18     pub sa_handler: usize,
19     /// 控制信号处理行为的标志位
20     pub sa_flags: SigActionFlag,
21     pub sa_restorer: usize,
22     /// 在执行信号处理函数期间临时阻塞的信号集合
23     /// - 信号处理函数执行时，内核会自动将 sa_mask 中的信号添加到进程的阻塞信号
24     集中
25     /// - 处理函数返回后，阻塞信号集恢复为原状态
26     pub sa_mask: SigMask,
27 }

```

代码 6-1 信号相关结构体定义

## 6.2 信号传输

在 Del0n1x 中，用户可以通过 kill 系统调用向进程传送信号，利用参数 pid 可以找到对应的进程或进程组，然后调用 TCB 成员函数接口 proc\_recv\_siginfo 将信号推入对应进程的待处理信号队列中（结构如下）：

```

1  pub struct SigPending {
2      /// 检测哪些 sig 已经在队列中,避免重复加入队列
3      mask: SigMask,
4      /// 普通队列
5      fifo: VecDeque<SigInfo>,
6      /// 存放 SIGSEGV, SIGBUS, SIGILL, SIGTRAP, SIGFPE, SIGSYS
7      prio: VecDeque<SigInfo>,
8      /// 如果遇到的信号也在 need_wake 中，那就唤醒 task
9      pub need_wake: SigMask,
10 }

```

代码 6-2 SigPending 结构体

我们将信号处理队列分为普通队列和优先队列,对不同的信号做了优先级处理,这样的数据结构时的 Del0n1x 对于紧急时间和高优先级时的相应延迟更低,提高了内核的实时性。

对于队列中的信号结构体 SigInfo 设计,我们借鉴了 Linux 中的 siginfo\_t 实现方式,同时对其进行了简化和封装,能够携带更多的数据信息(发送者 pid、子进程 exit code 和信号编码等),这样极大的方便了 Wait4 和 do\_signal 中对于不同信号的处理分发流程。

### 6.3 信号处理

进程因系统调用、中断或异常进入内核态,完成内核任务后,在返回用户态前,内核会检查该进程的未决信号。Del0n1x 中信号处理集中在 do\_signal 函数中,我们会依次遍历 prio 和 fifo 队列,如果该信号没有被阻塞,则根据 siginfo 中信号编码找到对应的信号处理函数 KSigAction,然后对 KSigAction 中的 sa\_type 字段进行模式匹配,对应的动作分别为 Ignore(忽略该信号)、Default(系统默认处理)和 Customized(用户自定义处理函数)。

对于用户自定义函数,内核会下面的流程进行处理(如下图):

构建用户态栈帧:在内核栈中创建新栈帧,如果用户没有自定义栈帧位置,那么默认为将用户栈 sp 向低地址扩展分配,确保信号处理函数有独立栈空间。

修改返回上下文:将原用户态执行点(如 pc 寄存器)保存到 UContext 中,然后复制到用户栈;然后修改当前 trap\_context 指向用户处理函数。

切换至用户态:跳转至信号处理函数入口并执行用户自定义函数。在这一过程中,为了避免信号的嵌套处理,需要将原信号加入屏蔽字。

切换到内核态:处理函数结束后调用 sigreturn() 系统调用,主动陷入内核态。内核从用户栈恢复原进程上下文,清除信号屏蔽字。

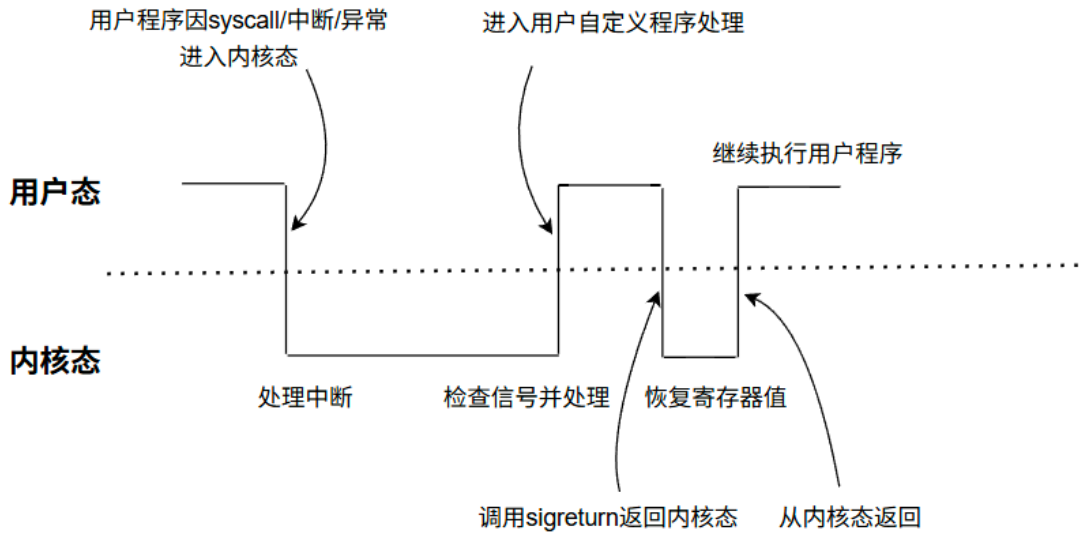


图 6-1 信号处理

## 第 7 章 时钟模块

### 7.1 定时器队列

在操作系统的定时器模块中，我们创造性地实现了一套混合定时器管理系统，将时间轮算法与最小堆优化相结合，相较于 Phoenix 和 Pantheon 使用的最小堆数据结构，我们的混合定时器在时间的处理上更加精细，同时在特定场景下更加的高效。

我们将定时器分为了短期和长期两种，以 600ms 作为阈值，将小于阈值的划分为短期定时器，大于阈值的划分为长期定时器。两种定时器分别存储在时间轮和最小堆中。

```
1 // 定时器队列
2 pub struct TimerQueue {
3     /// 短期定时器 (<600ms)
4     wheel: SpinNoIrqLock<TimingWheel>,
5     /// 长期定时器 (最小堆)
6     long_term: SpinNoIrqLock<BinaryHeap<TimerEntry>>,
7     /// 定时器句柄计数器
8     handle_counter: SpinNoIrqLock<u64>,
9 }
```

代码 7-1 TimerQueue 结构

### 7.1.1 时间轮设计

时间轮是一种高效的定时器管理数据结构，特别适合处理大量短周期定时器。鉴于我们内核的时钟中断间隔为 10ms，所以将时间轮划分为 60 个槽位，同时时间轮的滴答间隔为 10ms，10ms 自动推进一槽，与硬件时钟中断完美同步。如下图所示，当时间指针指向的槽位为 1 时，代表这次推进将处理 1 槽位中所有的定时器。槽内采用平铺向量存储，插入/删除操作达到  $O(1)$  常数时间，相较于最小堆插入  $O(\log n)$ /删除  $O(\log n)$  更加高效。

```
1 struct TimingWheel {
2     // 时间轮的槽数组，每个槽存储一组定时器条目
3     slots: [Vec<TimerEntry>; TIME_WHEEL_SLOTS],
4     // 当前指向的槽索引，随着时间推移循环递增
5     current_slot: usize,
6     // 时间轮当前表示的时间点
7     current_time: Duration,
8 }
```

代码 7-2 TimingWheel 结构体

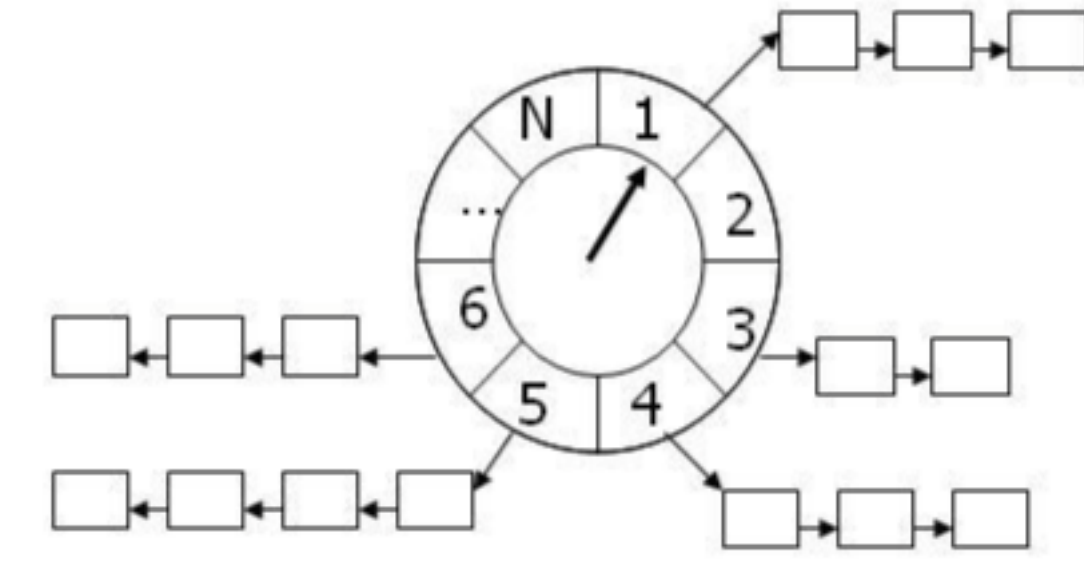


图 7-1 时间轮结构

时间轮推进算法如下所示，通过该函数我们获取到所有超时的定时器 waker，返回给上级调用者用于批量唤醒。

#### 算法 7-1: 时间轮推进算法

```
function advance_to(target_time: Duration) to Vec
```

---

```

1  wake_list ← 空列表
2  计算需要推进的槽数 slots_to_advance = calc_slot(target_time)
3  for 每个需要处理的槽 do
4      for 当前槽中的每个定时器条目 do
5          if 条目已过期(expire ≤ target_time) then
6              将条目的 waker 加入唤醒列表
7              从槽中移除该条目
8          else
9              保留该条目
10     移动指针到下一个槽
11     更新时间轮当前时间
12     if 当前时间超过目标时间 then 提前退出
13 return wake_list

```

---

## 7.2 定时器

我们使用 `TimeEntry` 表示定时器,每个定时器都携带专属的 `TimerHandle`——一个单调递增的唯一标识符。当异步 `Future` 结束其生命周期时,我们需要 `drop` 其中剩下时的定时器,这时就可以通过对比 `TimerHandle` 找到对应的定时器。系统先在时间轮中闪电扫描,然后扫描二叉堆。为优化堆内搜索,我们设计了临时缓存策略:将非目标项暂存后重新入堆,避免了重建整个堆的昂贵开销。这种双路径检索确保删除操作始终保持高效。

```

1  // 定时器条目
2  struct TimerEntry {
3      /// 到期时间
4      expire: Duration,
5      /// 唤醒器
6      waker: Option<Waker>,
7      /// 用于取消的句柄
8      handle: TimerHandle,
9  }

```

代码 7-3 `TimerEntry` 结构体

为了充分利用异步的优势,我们将需要监测的任务被封装在一个超时 `Future` 中(如下所示), `deadline` 表示任务的超时期限, `timer_handle` 用于 `Drop` 机制中确保定时器资源的回收,即使任务提前完成也不会留下幽灵定时器。我们使用 `poll` 轮循的方式对任务进行推进检测,轮询时执行三重检测:首先尝试推进内部任务,其次检查期限是否届满,最后才注册唤醒器,并且唤醒器只会注册一次。

```

1  pub struct TimeoutFuture<F: Future> {
2      inner: F,
3      deadline: Duration,
4      timer_handle: Option<TimerHandle>,
5  }
6
7  impl<F: Future> Drop for TimeoutFuture<F> {
8      fn drop(&mut self) {
9          // 确保定时器被删除
10         if let Some(handle) = self.timer_handle.take() {
11             TIMER_QUEUE.cancel(handle);
12         }
13     }
14 }

```

代码 7-4 TimeoutFuture 结构体与 Drop 实现

目前 Del0n1x 已实现较为高效的定时任务管理,但是该优化并没有在初赛测试用例中体现出来,在初赛测试用例中大部分定时器时间为 1s,其实都会被分配到最小堆结构中。这也反向说明了对于代码的优化也要在特定的场景中才能得到体现,比如时间轮适合处理大量短周期定时任务,最小堆适合处理少量长周期定时任务。

## 第 8 章 网络模块

在网络模块中,我们使用了官方推荐的 smoltcp 库作为网络协议栈的底层处理,该库具有十分高效的网络协议栈处理,同时通过 Rust 内存安全、事件驱动模型和高度模块化设计,成为资源受限内核场景的理想网络栈。

为了体现 linux 中"一切皆文件"的思想,我们为网络系统中的关键结构体 Socket 实现了 Filetrait,统一了网络和文件系统,方便内核设计中对其进行高效管理。在此基础上,我们对网络模块实现分层架构设计,从下至上依次为物理网络驱动层、传输层、Socket 接口层。

目前在初赛阶段我们实现了一个简单的本地回环网络,但是由于时间原因,还未适配通过 netperf、iperf 网络测例,不过这也是我们在决赛阶段的一个目标,努力实现一个简洁高效、易于学习的网络系统模块。

### 8.1 Socket 套接字

系统核心围绕 SockMeta 元数据结构和 Socket trait 展开,通过实现 FileTrait 将网络套接字无缝集成到文件系统中,实现了 Linux"一切皆文件"的设计哲学。整个架

构体现了协议无关性设计，TCP 与 UDP 套接字通过统一的接口向上提供服务，底层差异由各自的具体实现处理。为了利用 sockfd 操控套接字，我们将 socket 映射到 fdtable 中：

```

1  /// 将一个 socket 加入到 fd 表中
2  pub fn sock_map_fd(socket: Arc<dyn FileTrait>, cloexec_enable: bool) -> SysResult<usize> {
3      let mut flag = OpenFlags::O_RDWR;
4      let fdInfo = FdInfo::new(socket, flag);
5      let new_info = fdInfo.off_Ocloexec(!cloexec_enable);
6      let task = current_task().expect("no current task");
7      let fd = task.alloc_fd(new_info)?;
8      Ok(fd)
9  }
```

代码 8-1 将 socket 映射到 fd 的函数

SockMeta 结构体作为套接字的核心元数据容器，采用 Rust 的强类型系统精确描述套接字状态。其字段设计反映了网络连接的全生命周期：从初始创建时的空端口和端点，到绑定后的本地地址确定，再到连接建立后的远程端点记录。特别是通过 Option 类型明确区分已初始化和未初始化状态，避免了传统 C 实现中常见的无效值问题。

```

1  pub struct SockMeta {
2      pub domain: Sock,
3      pub iptype: IpType,
4      pub recv_buf_size: usize,
5      pub send_buf_size: usize,
6      pub port: Option<u16>,
7      pub shuthow: Option<ShutHow>,
8      pub local_end: Option<IpEndpoint>,
9      pub remote_end: Option<IpEndpoint>,
10 }
```

代码 8-2 SockMeta 结构体

Socket trait 定义了完整的套接字操作接口，既包含标准的 bind/connect/listen 等基本操作，也提供了 send\_msg/recv\_msg 等增强功能。该 trait 标记为 async\_trait 以适应现代异步 IO 需求，同时继承 FileTrait 实现了文件描述符的统一管理，可以通过 Filetrait 中的 pollin、pollout 异步检查网络缓冲区是否可读可写。



```

1  #[async_trait]
2  pub trait Socket: FileTrait {
3      /// 异步接受一个传入的连接请求
4      async fn accept(&self, flags: OpenFlags) -> SysResult<(IpEndpoint, usize)>;
5      /// 异步连接到指定的地址
6      async fn connect(&self, addr: &SockAddr) -> SysResult<()>;
7      /// 异步发送消息到指定地址
8      async fn send_msg(&self, buf: &[u8], dest_addr: &SockAddr) -> SysResult<usize>;
9      /// 异步接收消息
10     async fn recv_msg(&self, buf: &mut [u8]) -> SysResult<(usize, SockAddr)>;
11     /// 绑定套接字到本地地址
12     fn bind(&self, addr: &SockAddr) -> SysResult<()>;
13     /// 开始监听传入连接
14     fn listen(&self, backlog: usize) -> SysResult<()>;
15     /// 设置接收缓冲区大小
16     fn set_recv_buf_size(&self, size: u32) -> SysResult<()>;
17     /// 设置发送缓冲区大小
18     fn set_send_buf_size(&self, size: u32) -> SysResult<()>;
19     /// 获取当前接收缓冲区大小
20     fn get_recv_buf_size(&self) -> SysResult<usize>;
21     /// 获取当前发送缓冲区大小
22     fn get_send_buf_size(&self) -> SysResult<usize>;
23     /// 关闭套接字
24     fn shutdown(&self, how: ShutHow) -> SysResult<()>;
25     /// 获取套接字绑定的本地地址
26     fn get_sockname(&self) -> SysResult<SockAddr>;
27     /// 获取对端连接的地址
28     fn get_peername(&self) -> SysResult<SockAddr>;
29     /// 设置 TCP 保持活动选项
30     fn set_keep_alive(&self, action: u32) -> SysResult<()>;
31     /// 启用/禁用 Nagle 算法(TCP 延迟发送)
32     fn enable_nagle(&self, action: u32) -> SysResult<()>;
33     /// 获取套接字类型
34     fn get_socktype(&self) -> SysResult<Sock>;
35 }

```

代码 8-3 Socket trait 定义

## 8.2 Ethernet 设备

Del0n1x 使用 NetDev 结构体实现网络设备的封装和抽象，其包含两个关键组件：device 字段表示具体的网络设备类型，当前支持环回接口；iface 字段维护了访问 smoltcp 协议栈的通道。这种分离设计使得设备驱动与协议栈保持松耦合，未来扩展新设备类型时无需修改上层协议逻辑。

```

1 pub enum NetDevType {
2     Loopback(Loopback),
3     Unspec,
4 }
5
6 pub struct NetDev {
7     pub device: NetDevType,
8     pub iface: Interface,
9 }

```

代码 8-4 NetDevType 与 NetDev 结构体

#### ◆ 网络轮循

poll()方法是整个网络栈的驱动引擎,实现了事件处理的核心循环,每次 poll 轮循将处理 SOCKET\_SET 中所有的句柄。

iface.poll()调用实现了三层重要功能:

- 接收处理:从设备读取数据包并递交给相应协议处理程序
- 发送处理:将协议栈待发送数据提交给设备驱动
- 状态更新:维护 TCP 定时器、重传队列等状态机

该方法主要在发送和就收数据的循环中使用,驱动协议栈与设备的异步交互。

```

1 pub fn poll(&mut self) {
2     let instant = Instant::from_millis(get_time_ms() as i64);
3     let mut socket = SOCKET_SET.lock();
4     let device = match self.device {
5         NetDevType::Loopback(ref mut dev) => dev,
6         NetDevType::Unspec => panic!("Device not initialized"),
7     };
8     self.iface.poll(instant, device, &mut socket);
9 }

```

代码 8-5 网络轮询 poll 方法

### 8.3 传输层——UDP 与 TCP

TCP 是一种面向连接的字节流套接字,而 UDP 是一种无连接的报文套接字。他们都继承了 SockMeta 中的字段,实现了不同的 Socket trait。

```

1 pub struct TcpSocket {
2     pub handle: SocketHandle,
3     pub flags: OpenFlags,
4     pub sockmeta: SpinNoIrqLock<SockMeta>,
5     pub state: SpinNoIrqLock<TcpState>,
6 }
7
8 pub struct UdpSocket {
9     pub handle: SocketHandle,
10    pub flags: OpenFlags,
11    pub sockmeta: SpinNoIrqLock<SockMeta>,
12 }

```

代码 8-6 TcpSocket 与 UdpSocket 结构体

在 Tcp、Udp 数据结构中，都是用 SocketHandle 表示自己的句柄。该句柄存放在 smoltcp 提供的全局 SOCKET\_SET 中，该 set 可以类比为 fdtable，方便快速获取到该套接字的相关句柄，并通过该句柄实现与 smoltcp 底层网络栈功能交互。

```

1 /// 全局 handle 管理器
2 pub static ref SOCKET_SET: SpinNoIrqLock<SocketSet<'static>> =
3     SpinNoIrqLock::new(SocketSet::new(vec![]));
4
5 /// 通过闭包快速获取到对应的 handle
6 pub fn with_socket<F, R>(&self, f: F) -> R
7 where
8     F: FnOnce(&mut udp::Socket<'_>) -> R,
9 {
10     let mut binding = SOCKET_SET.lock();
11     let socket = binding.get_mut::<udp::Socket>(self.handle);
12     f(socket)
13 }

```

代码 8-7 SOCKET\_SET 及 with\_socket 函数

## 8.4 Port 端口分配

在网络交互过程中，PortManager 负责动态管理 TCP/UDP 端口资源的分配与回收。该系统采用多层级管理策略，在保证线程安全的前提下实现高效端口分配。为了记录端口分配情况，Del0n1x 采用了双位图的设计，单个协议仅需 8KB 内存（65536 位）即可实现  $O(1)$  复杂度的状态查询，相比传统哈希表节省了 90% 以上的内存开销。双位图设计彻底避免了 TCP/UDP 端口冲突。

```

1  pub struct PortManager {
2      /// 动态端口范围
3      pub start: u16,
4      pub end: u16,
5      /// 回收端口队列
6      pub recycled: VecDeque<u16>,
7      /// TCP 端口位图
8      pub tcp_used_ports: BitVec,
9      /// UDP 端口位图
10     pub udp_used_ports: BitVec,
11 }

```

代码 8-8 PortManager 结构体

## ◆ 分配算法设计

alloc 方法实现了三级分配策略，如算法所示。先从回收队列获取，，用局部性原理提升缓存命中率；然后随机尝试，PORT\_RANGE 范围内进行有限次随机探测；如果都失败，最后使用顺序扫描，保在极端情况下仍能穷尽搜索。

**算法 8-1: 端口分配三级策略**

```

function alloc(domain: Sock) to Result
1  if recycled queue is not empty then
2      port ← recycled.pop_front()
3      mark_used(domain, port)
4      return port
5  chance ← (end - start) - |recycled|
6  for i = 0 to chance-1 do
7      random_port ← start + (random() % PORT_RANGE)
8      if try_mark_used(domain, random_port) then
9          return random_port
10 for port ← start to end do
11     if try_mark_used(domain, port) then
12         return port
13 return error EADDRINUSE

```

## 第 9 章 设备

设备的管理在操作系统中起到至关重要的部分。将设备进行适当的抽象和封装，确保系统对设备进行规范的管理，并为其他模块提供简便易用的接口，并为拓展新的设备提供方便。

设备管理模块包括设备的发现和初始化、驱动程序匹配与加载等功能。操作系统需要设备管理系统与计算机系统其他部分进行信息的交换，直接影响整个系统的可用性、稳定性与性能。

在我们的系统中，基本设备抽象为基本设备属性 BaseDriver trait，所有的设备均要实现这个接口。

并拓展出块设备、tty 设备、网络设备等

例如块设备抽象为块设备属性 BlockDriver trait，这个属性继承了基本设备属性 BaseDevice trait

对于一个块设备的具体实现，需要实现块设备属性 BlockDriver trait，通过这样的设计，可以将设备进行一个统一的管理与使用。

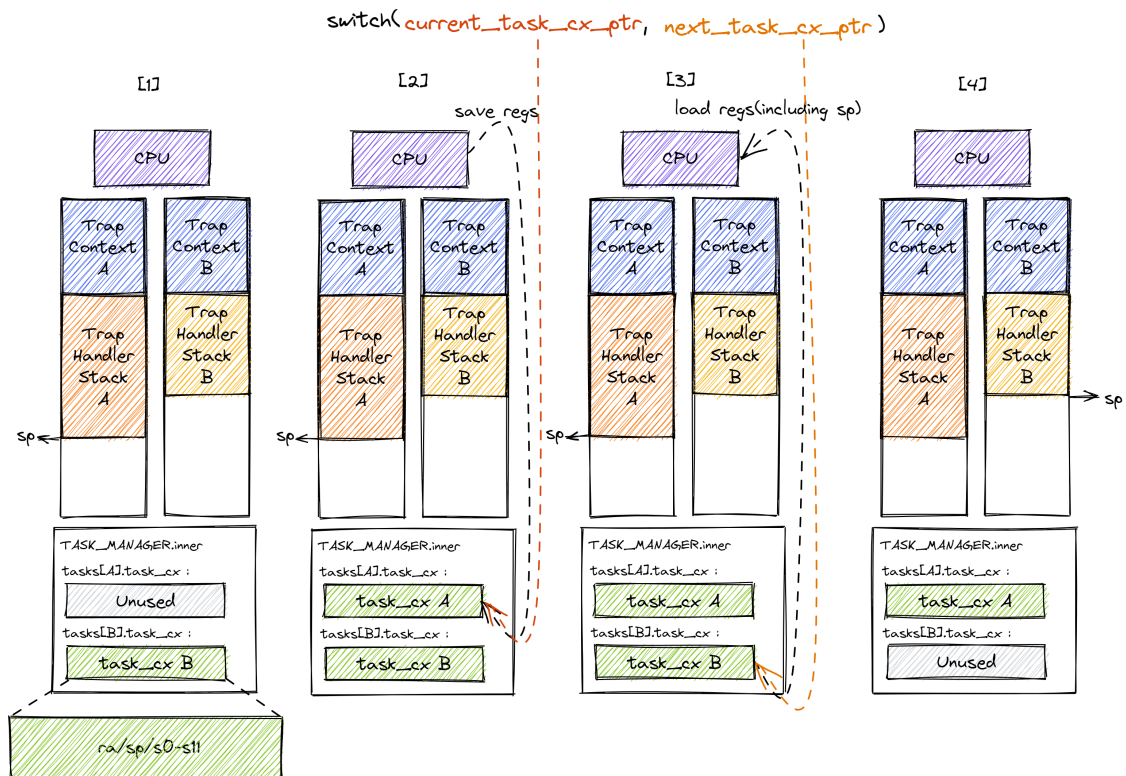


图 9-1 rCore 上下文切换

## 9.1 设备树

操作系统通过设备树实现对设备的发现

操作系统内核获取到设备树的地址的流程如下：

1. OpenSBI 启动：当系统启动时，OpenSBI 固件首先运行。它完成基础的硬件初始化，如内存控制器设置、I/O 初始化等。
2. 传递控制权到内核：OpenSBI 初始化完成后，将控制权传递给内核的入口点，并传递必要的参数。这些参数包括：
  - `hart_id`：当前硬件线程的 ID。
  - `dtb_addr`：设备树地址，该地址指向设备树描述符 (DTB)，描述了系统的硬件布局和配置信息。

对于 loongarch qemu virt 平台，我们通过查看 qemu 源码确定了设备树的地址

通过 `probe` 方法，从设备树的根节点进行搜索，获得当前机器上的设备与其地址。当检测到系统支持的设备的时候就调用注册函数，使用对应的驱动程序完成对这一设备的抽象，形成内核中的设备实体，供给其他模块使用。

未来我们希望完善 `devtmpfs` 实现从设备到文件的抽象。

通过设备树解析，我们可以实现同一份内核二进制在不同的硬件上启动。

## 第 10 章 硬件抽象层

### 10.1 硬件抽象层 (Hardware Abstraction Layer, HAL) 总览

为支持多平台运行与测试，Del0n1x 在内核中实现了功能基本完整的硬件抽象层，可以运行于 RISC-V64 和 LoongArch64 两种指令集架构的 QEMU 平台。Del0n1x 的内核代码全部基于硬件抽象层开发，屏蔽了架构细节和平台差异，具有更好的兼容性和可移植性。

【硬件抽象层示意图】

### 10.2 处理器访问接口

Del0n1x 在 `riscv` 和 `loongarch` 外部库的帮助下，对 RISC-V 和 LoongArch 两种架构下的处理器访问实现了统一抽象。在内核代码中，对处理器及相关资源的访问包括以下几种类型：

- 读写通用寄存器
- 访问控制状态寄存器
- 执行特殊的控制指令（如内核刷表指令）

在 Del0n1x 中, 这些操作只需要调用统一的接口即可完成, 具体实现细节位于 HAL 中, 便于构建架构无关的内核代码。

### 10.3 内核入口例程

RISC-V 和 LoongArch 架构在内核启动上的细节有所异同。LoongArch 架构启动时需要配置直接映射地址翻译模式, 需要为分页地址翻译模式设置高半空间、低半空间两个页表 token, 还需要设置 TLB 重填异常入口, 这些都是 RISC-V 架构下内核启动所不需要的步骤。Del0n1x 的硬件抽象层将初始化时的架构相关部分抽象为 arch\_init() 例程, 供内核初始化时调用, 为内核入口的规范化提供支持。【代码, arch\_init 重构完毕后填写】

### 10.4 内存管理单元与地址空间

#### 10.4.1 物理内存

RISC-V 和 LoongArch 的 QEMU virt 平台的物理编址方式和布局存在很大区别。以下分别为两个平台的物理地址布局：【图片】

RISC-V QEMU virt 平台将 RAM 编码于 0x8000\_0000 以上的物理地址空间。而 LoongArch QEMU virt 平台则将 RAM 切分为 lowram 和 highram 两个部分, 其中 lowram 位于 0x1000\_0000 以下的物理地址空间, highram 位于 0x9000\_0000 以上的物理地址空间。Del0n1x 将内核镜像加载于 lowram 中。

#### 10.4.2 分页地址翻译模式

Del0n1x 使用 SV39 分页地址翻译模式。在不考虑直接映射窗口的前提下, Del0n1x 将完整的 39 位虚拟地址空间分为高半部分和低半部分, 高半部分的地址范围为 0xffff\_ffc0\_0000\_0000~0xffff\_fff\_fff\_fff, 低半部分的地址范围为 0x0000\_0000\_0000\_0000~0x0000\_003f\_fff\_fff。【可以用表格展示】

RISC-V 架构使用控制状态寄存器 SATP 控制分页模式类型 (如 SV39) 并存储页表根目录的物理页号。Del0n1x 为每一个进程创建一张页表, 并为其应用整个虚拟地址空间内的全部映射。当地址空间切换时, 通过对 SATP 的修改, 即可达到切换页表的目的。

LoongArch 架构的情况有所不同。LoongArch 架构使用 CSR.CRMD 控制状态寄存器的 PG 位开关分页地址翻译模式, 通过设置 CSR.PWCL 和 CSR.PWCH 控制状态寄存器中的 Dir{}\_base、Dir{}\_width 等部分手动配置分页模式细节。同时, LoongArch

架构下虚拟地址空间低半部分和高半部分的遍历应用不同的页表，其根目录地址分别经由 CSR.PGDL 和 CSR.PGDH 控制状态寄存器设置。这意味着 LoongArch 架构下一个进程需要创建两张页表，分别映射虚拟地址空间的两半。当地址空间切换时，需要按需修改 CSR.PGDL 或 CSR.PGDH，从而灵活保留不需要切换的部分。

Del0n1x 在硬件抽象层中为两种架构实现了分页地址翻译模式的初始化，为 LoongArch 手动配置了 SV39 分页模式。由于 Del0n1x 使用虚拟地址空间的高半部分作为通用的内核地址空间，借助 LoongArch 架构的“双目录”设计，在地址空间切换时只需修改 CSR.PGDL。

【代码：mmu\_init】

#### 10.4.3 页表

Del0n1x 支持两种架构下的多级页表。在硬件抽象层中，Del0n1x 为两种架构分别实现了对页表项（Page Table Entry，PTE）的封装，可以便捷地访问 PTE 中存储的物理页号和各标志位。Del0n1x 使用 Rust 宏为 PTE 标志位定义了统一的 checker 和 setter 方法，用于在内核代码中灵活修改 PTE 标志位。

Del0n1x 的设计中，LoongArch 架构的内核地址空间需要单独使用一张页表，而 RISC-V 架构的内核地址空间则与用户地址空间共用同一张。Del0n1x 在硬件抽象层中为两个架构均映射一张内核页表，区别在于 LoongArch 架构下该页表会被写入 CSR.PGDH 并永不切换，而 RISC-V 架构下该页表只作为一个映射用户地址空间时的模板。当创建新的进程页表时，内核页表中根目录中映射地址空间高半部分的页表项将被复制到新的页表中，供内核态访问。

【代码，page\_table】

#### 10.4.4 直接映射窗口

LoongArch 架构支持直接映射地址翻译模式，该模式下允许通过修改 CSR.DMW0 ~ CSR.DMW3 控制状态寄存器配置至多 4 个直接映射窗口。当虚拟地址的高 4 位（在 LoongArch64 中为[63:60]位）恰好与某个直接映射窗口的高 4 位相同时，虚拟地址将被直接映射为其低 PALEN 位的物理地址（PALEN 为机器物理地址长度）。使用 cpucfg 指令可以查明，LoongArch QEMU virt 平台下 PALEN=48，故直接映射窗口将 0xW000\_XXXX\_XXXX\_XXXX 范围内的所有虚拟地址映射为 0x0000\_XXXX\_XXXX\_XXXX。通过修改 CSR.DMWx 还可配置直接映射窗口的允许访问特权级、存储访问类型。



Del0n1x 使用 0x8000\_XXXX\_XXXX\_XXXX 和 0x9000\_XXXX\_XXXX\_XXXX 两个直接映射窗口，均限定内核特权级（PLV0）使用，分别用于设备访问和物理内存访问。

#### 10.4.5 TLB 重填

LoongArch 架构使用软件管理 TLB。当发生 TLB 中没有匹配项时，将触发 TLB 重填异常，跳转到内核设置的 TLB 重填入口执行软件重填。Del0n1x 使用了往届优秀作品 NPUCore-IMPACT 编写的 TLB 重填代码。

【可能贴代码】

## 第 11 章 总结与展望

### 11.1 工作总结

1. 支持 loongarch64 和 riscv64 架构，实现了自己的硬件抽象层 hal。
2. 实现进程管理，以无栈协程的方式高效调度任务。
3. 实现设备树解析
4. 实现 COW、懒分配等内存优化策略。
5. 实现虚拟文件系统，将具体文件系统与内核解耦合，实现了 pagecache 和 dentry cache 优化。
6. 支持信号模块，实现进程间信号通信。
7. 时钟模块支持时间轮混合最小堆的时钟管理方式。
8. 支持多核调度运行。
9. 实现 100+ 调系统调用。

### 11.2 未来计划

1. 完善 net 模块，支持网络上板。
2. 完善 loop 设备，实现功能更加完善的 mount 机制。
3. 适配龙芯板和 riscv 板，完善相关驱动。
4. 支持外设中断。
5. 支持更多 ltp 测例，修复更多内核不稳定的 bug。
6. 支持更多现实应用。

## 11.3 参考

- Phoenix：无栈协程、内存管理
- polyhal：硬件抽象层