



# NoAxiom OS

*Operating System Design Manual*

陈润澎, 刘文涛

August 17, 2025

# Contents

<b>一、概述</b>	<b>5</b>
1.1 系统简介 . . . . .	5
1.2 系统整体架构 . . . . .	5
1.3 系统完成情况 . . . . .	7
<b>二、任务调度</b>	<b>9</b>
2.1 无栈协程特性 . . . . .	9
2.1.1 Rust 的异步语法支持 . . . . .	9
2.1.2 异步运行时 . . . . .	10
2.2 模块化调度抽象 . . . . .	11
2.2.1 调度器特性抽象 . . . . .	11
2.2.2 运行时特性抽象 . . . . .	12
2.3 多种调度器设计 . . . . .	12
2.3.1 先入先出调度 . . . . .	12
2.3.2 双优先级调度 . . . . .	13
2.3.3 分步队列调度 . . . . .	14
2.3.4 多级调度 . . . . .	14
2.4 实时调度性能调优 . . . . .	15
2.5 调度信息维护 . . . . .	17
2.5.1 调度实体 . . . . .	17
2.5.2 调度优先级 . . . . .	17
2.6 异步运行时实现 . . . . .	18
2.6.1 异步任务的生成 . . . . .	18
2.6.2 异步任务的调度与运行 . . . . .	18
<b>三、进程管理</b>	<b>20</b>
3.1 任务控制块 . . . . .	20
3.1.1 资源访问属性 . . . . .	20
3.1.2 任务控制块 . . . . .	20
3.2 任务状态 . . . . .	23
3.2.1 任务状态转化关系 . . . . .	23
3.2.2 调度相关状态简化设计 . . . . .	23
3.3 任务间关系 . . . . .	24
3.3.1 进程间关系 . . . . .	24
3.3.2 线程间关系 . . . . .	24
3.3.3 任务管理器 . . . . .	25

3.3.4 进程组管理器	25
3.4 任务生命周期	25
3.5 任务的创建	26
3.5.1 初始进程的加载	26
3.5.2 普通任务的加载	27
3.6 任务运行流程	28
3.6.1 任务切换	28
3.6.2 任务主函数	29
3.7 任务资源释放	30
<b>四、内存管理</b>	<b>32</b>
4.1 地址空间结构	32
4.1.1 内核地址空间初始化	32
4.1.2 内核地址空间排布	32
4.1.3 用户地址空间排布	33
4.2 内存管理器	33
4.2.1 堆分配管理器	33
4.2.2 物理页帧管理器	34
4.2.3 资源占用维护	35
4.3 地址空间维护	35
4.4 延迟分配技术	36
4.4.1 写时复制	36
4.4.2 懒分配	37
4.4.3 合法性检查与分配	37
4.5 用户指针	38
4.5.1 基于异常的指针预检查	38
4.5.2 用户指针的延迟分配与异步读取	41
<b>五、指令集架构抽象层</b>	<b>43</b>
5.1 双架构支持概览	43
5.2 硬件架构解耦设计	43
5.3 硬件初始化	44
5.3.1 启动时初始化	44
5.3.2 内核功能初始化	46
5.4 页表管理	46
5.4.1 页表抽象	46
5.4.2 页表项抽象	47
5.4.3 TLB 维护	50
5.5 异常处理	51

5.5.1 异常处理抽象 . . . . .	51
5.5.2 异常类型抽象 . . . . .	51
5.5.3 异常类型转化 . . . . .	51
5.6 上下文维护 . . . . .	52
5.6.1 寄存器抽象 . . . . .	52
5.6.2 上下文抽象 . . . . .	53
<b>六、设备抽象层</b>	<b>55</b>
6.1 设备抽象设计 . . . . .	55
6.1.1 设备抽象层概览 . . . . .	55
6.1.2 设备类型抽象 . . . . .	55
6.1.3 设备驱动抽象 . . . . .	56
6.2 设备探测 . . . . .	56
6.2.1 设备树探测 . . . . .	56
6.2.2 总线设备嗅探 . . . . .	58
6.2.3 PCI 总线下的设备嗅探 . . . . .	58
6.2.4 MMIO 总线下的设备嗅探 . . . . .	59
6.3 驱动层 . . . . .	60
6.3.1 驱动层结构 . . . . .	60
6.4 异步块设备驱动 . . . . .	61
6.5 PLIC 中断控制器 . . . . .	63
6.5.1 初始化 . . . . .	63
6.5.2 外部中断处理 . . . . .	64
6.6 LoopBack 回环设备 . . . . .	65
<b>七、文件系统</b>	<b>67</b>
7.1 VFS . . . . .	67
7.1.1 总体结构设计 . . . . .	67
7.1.2 目录项 Dentry . . . . .	67
7.1.3 元信息 Inode . . . . .	71
7.1.4 文件 File . . . . .	73
7.1.5 超级块 SuperBlock . . . . .	74
7.2 缓存设计 . . . . .	75
7.2.1 块缓存 . . . . .	75
7.2.2 页缓存 . . . . .	76
7.2.3 路径缓存 . . . . .	83
7.2.4 性能对比 . . . . .	83
7.3 文件系统实例 . . . . .	88
7.3.1 具体文件系统 EXT4/FAT32 . . . . .	88

7.3.2 管道文件 . . . . .	89
7.3.3 套接字文件 . . . . .	91
7.3.4 基于内存的文件系统 RamFs . . . . .	93
7.4 虚拟文件 . . . . .	93
7.4.1 进程管理 /proc . . . . .	93
7.4.2 虚拟设备 /dev . . . . .	94
7.5 系统调用实现 . . . . .	95
7.5.1 I/O 读写类 . . . . .	96
7.5.2 I/O 多路复用 . . . . .	96
7.6 其他相关结构 . . . . .	97
7.6.1 文件描述符表 FdTable . . . . .	97
7.7 文件一致性管理 . . . . .	98
7.8 文件权限管理 . . . . .	98
7.8.1 权限结构设计 . . . . .	99
7.8.2 权限检查功能 . . . . .	99
<b>八、信号系统</b>	<b>102</b>
8.1 信号管理 . . . . .	102
8.2 信号处理流程 . . . . .	102
8.3 可中断系统调用 . . . . .	103
<b>九、网络模块</b>	<b>105</b>
9.1 结构概述 . . . . .	105
9.2 套接字池 SOCKET_SET . . . . .	108
9.2.1 套接字状态更新 . . . . .	108
9.2.2 套接字生命周期 . . . . .	109
9.3 端口管理器 PortManager . . . . .	109
9.3.1 端口生命周期 . . . . .	110
9.3.2 端口复用 . . . . .	111
9.4 TcpSocket . . . . .	111
9.5 UdpSocket . . . . .	113
9.6 git 测试 . . . . .	114
<b>十、物理板卡支持</b>	<b>116</b>
10.1 LS2K1000 地址未对齐异常 . . . . .	116
10.2 多核启动问题 . . . . .	117
10.3 块设备驱动支持 . . . . .	118
<b>十一参考文献</b>	<b>120</b>

# 一、概述

## 1.1 系统简介

NoAxiom 操作系统<sup>1</sup>是一款由杭州电子科技大学 NoAxiom 团队<sup>2</sup>完全自主开发的基于 Rust 语言的宏内核操作系统，可同时于 RISC-V64 的 qemu/visionfive 平台和 LoongArch64 的 qemu/2k1000 平台运行。本操作系统自主实现了**进程管理**、**内存管理**、**文件系统**、**任务调度**、**信号系统**、**时钟模块**、**网络模块**、**硬件抽象层**等子模块。NoAxiom 基于 Rust 无栈协程异步语法实现**异步调度**，在 I/O 等方面取得了优异的性能水平。此外，NoAxiom 自主设计实现了完整的硬件抽象层，实现内核与底层硬件的完全解耦。

## 1.2 系统整体架构

NoAxiom 操作系统可分为四个层次：**机器层**、**硬件抽象层**、**内核实现层**、**用户层**，此外还有全局设置子模块用于在各个层级之间传递设置信息与常量定义信息。图 1 中展示了 NoAxiom 的整体架构设计图。

更具体地，NoAxiom 自底向上地将其架构分为如下四个层次：

### ➢ 机器层 / Machine Layer

我们的目标是构建一个能够运行于多平台、多架构上的操作系统。我们的内核需要运行于 RISC-V64 架构和 LoongArch64 架构，且需要考虑到同时运行于 QEMU 虚拟机平台与物理上板平台。这些平台与架构均被归于机器层（Machine Layer），将会在上层的硬件抽象层中被实现为统一的接口提供给内核子系统。

### ➢ 硬件抽象层 / Hardware Abstraction Layer (HAL)

硬件抽象层用于统一机器层的复杂环境，并向上层内核提供统一的功能接口。内核子系统通过调用硬件抽象层内部的多个子模块，实现与底层机器环境的完全解耦。这一设计在大赛要求双架构运行的背景下发挥了重要作用，使得我们能够完全脱离底层架构进行，一定程度上减少了处理内核在不同架构/平台下运行带来的额外编码工作量。具体而言，硬件抽象层又细分为指令集架构抽象层 `arch`，平台抽象层 `platform`，内存抽象层 `memory`，驱动抽象层 `driver` 等子模块，下面是详细描述。

1. **指令集架构抽象层**：我们基于对 LoongArch 和 RISC-V 架构的了解，将所有受指令集架构影响的函数调用、类型定义、常量定义均视为同一对象 `Arch` 的 `trait` 子类，并将所有架构相关的实现放置于 `lib/arch` 库当中提供给所有上层子模块使用。借此内核实现了指令集架构上的解耦。
2. **平台层**：为了同时运行于虚拟机平台与物理上板环境，我们构建了平台抽象层用于进行平台与内核的解耦。具体而言，平台抽象层会根据编译时给定的平台与架构，向上提供统一的常量定义接口，从而实现内核与运行平台的完全解耦。
3. **内存管理层**：内存初始化布局会受到硬件平台的影响，同时驱动抽象层也需要依赖于物理页帧分配、虚实地址等功能，因此我们创建了内存抽象层向上提供统一的访存接口。本层次仅用于平台相关的内存初始化，以及架构、平台相关的内存抽象，真正的内存管理仍然由上层内核维护。

<sup>1</sup><https://github.com/NoAxiom/NoAxiom-OS>

<sup>2</sup><https://github.com/NoAxiom>

4. **设备抽象层**: 在操作系统启动时，驱动抽象层会进行**自动的设备嗅探**，读取设备树，并且自动嗅探挂载于不同总线上的设备信息，并且注册相应驱动，供上层内核使用。特别的，对于 RISC-V 架构，我们还会配置 PLIC(Platform-Level Interrupt Controller)，为需要中断的设备注册中断，并提供相应的中断处理例程。实现了内核与下层具体物理设备的解耦。

#### ➤ 内核实现层 / Kernel Implementation Layer

内核实现层基于下层硬件抽象层的统一接口，使用统一的逻辑实现各类型模块的具体功能，并向用户提供完整的系统调用支持。除了细分的子模块外，我们与大多数操作系统一样，创建了统一的系统调用接口，用于统合操作系统内核提供的各个子模块的具体功能，并通过相关系统规范向用户提供统一的系统调用服务。

具体而言，内核实现层细分为以下若干子模块：

1. **进程管理子模块**: 进程管理子模块通过提供进程控制块与若干类型的任务管理器，向其他子模块提供基本的进程信息维护功能。
2. **任务调度子模块**: 任务调度子模块借助 Rust 的无栈协程语法，维护了分时多任务的异步调度器，实现了任务执行的全流程异步维护。
3. **文件系统子模块**: 文件系统子模块通过统一的 VFS 接口，向下统合多种不同的文件系统类型，向上提供统一的文件访问行为抽象；此外，文件系统结合无栈协程异步让权机制，实现了**耗时行为的异步调度**，以及**异步的 I/O 多路复用**。
4. **信号子模块**: 信号子模块通过维护进程控制块中的信号信息，实现了不同进程之间的异步通信。
5. **内存管理子模块**: 内存管理子模块向其他模块提供统一的内存管理接口，并通过基于硬件抽象层提供的异常检查，用户提供精确的异常指针检查机制与懒分配让权机制。
6. **时间管理子模块**: 时间管理子模块通过设置若干时钟管理器，实现了可靠的时钟计时功能。
7. **网络子模块**: 网络子模块实现了 TCP 和 UDP 两种套接字 API，支持 IPv4 和 IPv6 地址，在无栈协程架构下，实现了高并发下的准确高效的网络传输。

#### ➤ 用户应用层 / User Application Layer

用户应用层基于内核提供的系统调用接口实现多种功能的用户程序。特别地，我们的内核内嵌了一个初始进程的 elf 文件用于进行用户程序的初始执行过程。赛事官方提供的测例均运行于这一层。

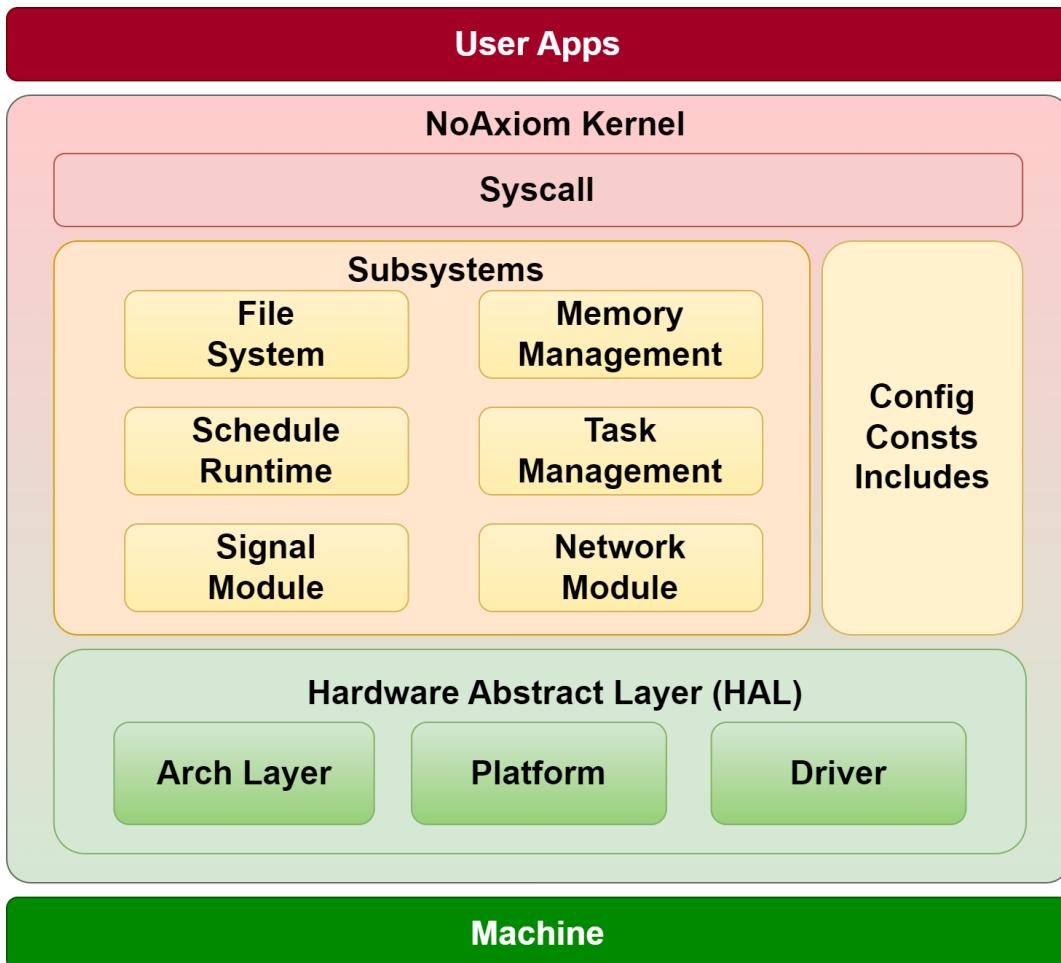


图 1: NoAxiom OS 整体架构

### 1.3 系统完成情况

截至目前，本系统共实现 112 条系统调用，涵盖了文件系统、I/O、网络、进程管理、信号处理、内存管理、调度管理、时间管理等各类系统调用类型。能够运行赛事官方提供的除 ltp 部分测试点以外的全部测例。

NoAxiom 在多个子模块均取得了较为明显的技术突破。截至目前，NoAxiom 的具体子模块完成情况如表 1 所示。

表 1: NoAxiom 系统子模块完成情况

子模块	实现情况
进程管理	统一的进程资源抽象
	<b>细粒度共享资源</b>
内存管理	内核与用户地址空间共享
	<b>懒分配与写时复制检查</b>
文件系统	用户指针合法性的快速检查
	<b>文件映射懒分配的完整异步让权</b>
任务调度	实现类 Linux 的 <b>VFS 虚拟文件系统</b>
	支持管道、套接字及各类虚拟文件的挂载
	支持 <b>异步 EXT4、FAT32 文件系统</b>
	为耗时读写操作实现了 <b>异步让权</b>
	实现了高效的 <b>页缓存加速文件读写</b>
信号系统	支持基于异步让权的 <b>I/O 多路复用</b>
	使用实现完整的 <b>分时多任务异步调度</b>
	使用统一的调度器特性抽象管理 <b>多种调度器</b>
架构抽象层	通过调度实体实现 <b>任务优先级</b>
	实现信号系统维护
设备抽象层	支持 <b>可中断系统调用的统一维护函数</b>
	<b>自主实现 RISC-V64 和 LoongArch64 双架构支持</b>
	<b>自主实现与内核统一解耦的顶层硬件抽象接口</b>
	<b>自主实现访存模块的架构解耦</b>
网络模块	<b>自主实现中断异常处理的架构解耦</b>
	<b>自主实现完备的设备管理抽象</b>
	实现不同架构下 <b>多总线的设备嗅探抽象管理</b>
	<b>自主实现基于外部中断的异步块设备驱动</b>
	为块设备驱动实现 <b>异步块缓存</b>
	支持 TCP 和 UDP 套接字
	支持 Ipv4 与 Ipv6 协议
	<b>实现端口复用</b>
	<b>支持等待时异步让权</b>

## 二、任务调度

### 2.1 无栈协程特性

Rust 语言于 2019 年推出了基于 `async / await` 关键字的无栈协程支持，简化了异步函数状态机的创建与维护过程，标志着 Rust 正式进入协程时代。而我们的内核也大量使用了这一语法特性，借助无栈协程实现了完整的异步让权机制。

#### 2.1.1 Rust 的异步语法支持

在 Rust 语言中，无栈协程底层由 `Future` 特性进行实现，其内部维护了一个无栈协程异步函数状态机，并通过调用 `Future::poll()` 方法进行异步函数的轮询操作。对应代码如下：

```
1 pub trait Future {
2     type Output;
3     fn poll(
4         self: Pin<&mut Self>, cx: &mut Context<'_>
5     ) -> Poll<Self::Output>;
6 }
```

上述代码中涉及若干关键对象，下面是它们的具体含义及解释：

1. `trait Future`：它是无栈协程中所有异步操作的上层抽象，当调用 `Future::poll()` 的时候，当前异步操作会尝试推动其进展，并返回其执行状态。
2. `enum Poll<T>`：异步操作的返回值，用于指示当前异步操作的进展。
  - `Poll::Pending`：表示当前异步操作仍未准备完毕，需要再次轮询。
  - `Poll::Ready(val)`：表示当前异步操作已经完成，并会附带准备完毕的返回数据。
3. `Pin<Self>`：`poll()` 方法中将当前对象包裹在 `Pin` 指针当中，该指针确保了其指向的对象位于固定的内存位置，不会被移动，这主要用于防止 `Future` 当中自引用导致的指针错误。
4. `Context`：其内部支持了上下文相关信息，如唤醒器 `Waker` 等信息的保存。它允许 `Future` 在返回 `Pending` 时在运行时内保存间接指向其执行上下文的指针，从而在反应器感知到对应事件时，唤醒挂起的任务并将其重新加入到执行器中。

而 Rust 中的 `async / await` 关键字是 `Future` 的语法糖，当一个函数被标记为 `async` 的时候，编译器会将其视作一个 `Future` 对象，并按照其内部的 `.await` 关键字排布，将其拆分为若干个子 `Future` 对象。

在执行异步函数时，程序会在堆空间上分配当前 `Future` 所需的数据空间，从而创建出无栈协程状态机。当执行器对于 `async` 函数生成的 `Future` 对象进行轮询的时候，其内部会按照当前状态机的执行状态，对当前执行到的子 `Future` 对象进行 `poll` 方法的调用，从而实现异步函数执行上下文的维护。

### 2.1.2 异步运行时

在 Rust 语言中，异步运行时（Async Runtime）由执行器（Executor），反应器（Reactor），调度器（Scheduler）等组件组成，这些组件共同构成了无栈协程运行时（Rust Stackless Coroutine Runtime）。我们的操作系统完全基于 Rust 无栈协程语法构建，所有任务最终均在运行时上完成执行。

下面对无栈协程运行时内部组件进行详细的描述。

- **执行器（Executor）**：执行器组件用于提供异步函数在同步控制流中的**执行功能**，当执行器执行一个任务时，任务将处于**执行态**。

Rust 的 `Future` 是惰性的，在调用其 `poll()` 方法前，它不会真正被驱动执行。Rust 语言规定了 `.await` 只能在在一个异步函数 `async fn` 内被调用，因此对于最外层的异步函数，需要一个执行器用于在同步控制流当中启动该无栈协程异步函数。在 NoAxiom 当中，我们选择了基于 `async_task`<sup>3</sup> 外部库构建无栈协程执行器。该外部库会通过 `AtomicUsize` 原子访存对于任务的调度状态进行内部维护，保证任务对象在重复唤醒时的唯一性。

- **反应器（Reactor）**：反应器用于保存异步执行时**挂起**的任务，并在对应事件完成时重新唤醒对应任务。当任务被保存在反应器内时，它将处于**挂起状态**。

在一个设计完备的异步系统当中，当一个 `Future` 对应的异步执行任务尚未完成时，系统应当将当前任务切出，并将任务保存在反应器当中。反应器会通过中断或轮询等方式监听对应的事件，当检测到对应的异步任务被标记为执行完成时，再重新唤醒对应的任务使其继续执行。

在 NoAxiom 当中，反应器并不对应一个特定组件，而是通过中断处理例程、时间管理例程等形式分布在内核中的各个组件内部。

- **调度器（Scheduler）**：调度器用于保存已被唤醒、**可运行**的任务。它可以向执行器提供可执行任务，也可以响应反应器唤醒的任务。当任务位于调度器内时，它将处于**可执行态**。

在部分 Rust 无栈协程设计当中，对于执行器和调度器并没有进行良好区分。但为了设计解耦，NoAxiom 选择将调度器与执行器严格区分以独立解耦内核调度过程，并构建了良好的调度器抽象，以实现多种任务调度功能。

异步运行时的内部组件构成了从调度器到执行器，执行器到反应器，反应器再唤醒回到调度器的完整执行循环。其执行循环如图 2 所示，其中描述了各个组件及对应的任务执行状态，可以参考 3.2 中的任务状态切换查看其具体对应关系。

---

<sup>3</sup><https://github.com/smoltc/smoltc/tree/main/src/task>

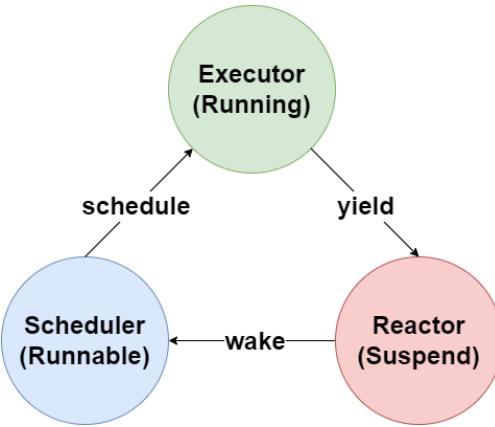


图 2: 异步运行时组件间关系

## 2.2 模块化调度抽象

我们的操作系统基于轻量级异步执行器开源库async-task构建了独属于我们操作系统的模块化调度器，实现了良好的调度器抽象，并针对不同调度场景设计了特定的调度器进行灵活调度。

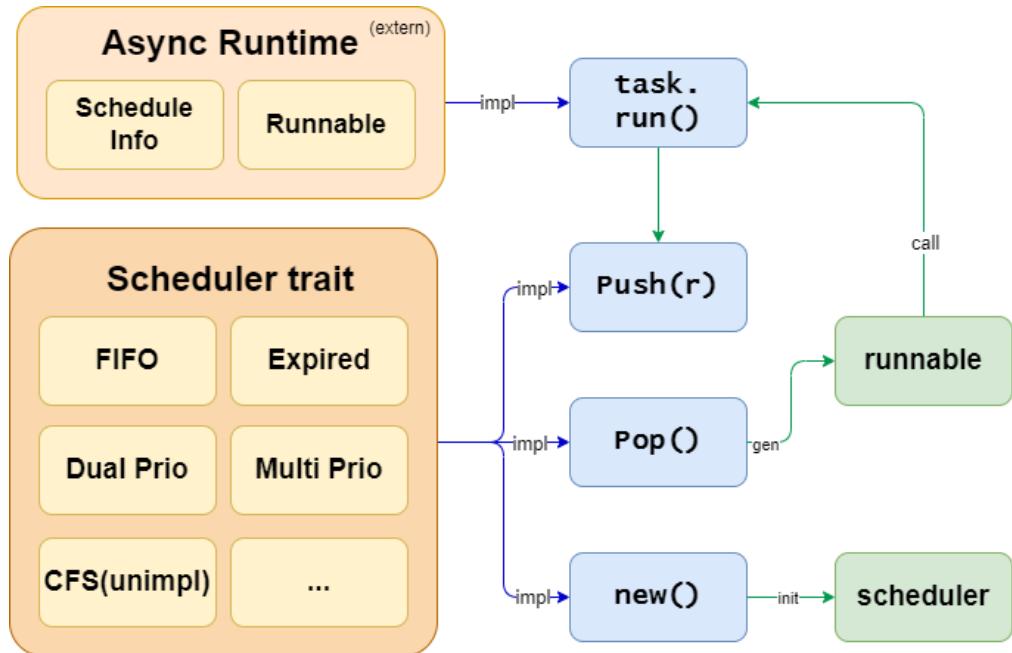


图 3: 模块化调度器设计

### 2.2.1 调度器特性抽象

根据我们在 2.1.2 中对于异步运行时的设计，在 NoAxiom 中，调度器最终被抽象为独立、统一的特性对象 `trait Scheduler<R>`。除了用于初始化的 `Scheduler::new()`<sup>4</sup>之外，`Scheduler` 中定义了两个关

<sup>4</sup>由于实现较为简单，下文中我们将省略展示的代码中对该方法的实现

键方法 `push()` 和 `pop()`，分别用于反应器中的 `Waker` 唤醒与执行器中的可执行任务获取过程。

具体函数定义如下所示。

```
1 // NoAxiom/kernel/src/sched/vsched.rs
2 pub trait Scheduler<R> {
3     fn new() -> Self;
4     fn push(&mut self, runnable: Runnable<R>, info: ScheduleInfo);
5     fn pop(&mut self) -> Option<Runnable<R>>;
6 }
```

基于该调度器特性抽象，我们可以设计出多种适用于不同调度场景的调度器，并向外暴露统一的接口给运行时以进行调度。

### 2.2.2 运行时特性抽象

在 `NoAxiom` 的设计当中，运行时 `Runtime` 用于统合执行器、调度器、反应器三者的功能，并向内部组件提供顶层的接口支持。其中具体成员方法及其含义如下：

- `run()` 方法：用于调用执行器进行单个任务的执行。
- `schedule()` 方法：用于反应器访问调度器插入可执行任务。
- `spawn()` 方法：用于将异步函数与任务控制块打包为统一的 `Runnable` 对象投入调度中。

```
1 pub trait Runtime<T, R>
2 where
3     T: Scheduler<R>,
4 {
5     fn new() -> Self;
6     fn run(&self);
7     fn schedule(&self, runnable: Runnable<R>, info: ScheduleInfo);
8     fn spawn<F>(&self: &'static Self, future: F, task: Option<&Arc<Task>>)
9     where F: Future<Output: Send + 'static> + Send + 'static;
10 }
```

通过提供运行时特性抽象，我们可以根据实际运行的平台选定特定的运行策略，并为运行时选定具体的底层调度器，以实现灵活的调度策略选择。

## 2.3 多种调度器设计

基于 `Scheduler<R>` 特性，我们可以方便的设计出多种底层调度器实现。下面将列举 `NoAxiom` 系统中的具体调度器实现。

### 2.3.1 先入先出调度

先入先出调度（First In First Out, FIFO）是一种常见且朴素的调度策略，通过对于所有任务的轮询过程避免了因调度优先级导致的任务饥饿情况。在 `NoAxiom` 当中，我们选择使用全局自旋锁 + 双端队列维护 FIFO 调度队列。具体代码如下。

```

1 // NoAxiom/kernel/src/sched/scheduler.rs
2 struct FifoScheduler {
3     queue: VecDeque<Runnable<Info>>,
4 }
5 impl Scheduler<Info> for FifoScheduler {
6     fn push(&mut self, runnable: Runnable<Info>, _: async_task::ScheduleInfo) {
7         self.queue.push_back(runnable);
8     }
9     fn pop(&mut self) -> Option<Runnable<Info>> {
10         self.queue.pop_front()
11     }
12 }

```

### 2.3.2 双优先级调度

在常规的调度过程中，有两类调度行为较为常见：

- 让权 (yield): 主要用于轮询行为，程序通过调用 `yield_now().await` 主动让出当前任务执行权，此时的任务应当被重新加入到调度队列末尾，等待后续任务执行完毕后再次进行轮询。
- 挂起-唤醒 (suspend-wake): 主要发生在异步操作的执行过程当中，如块设备访问、管道 I/O 等。该调度行为在唤醒时应当得到尽可能早的响应，以降低系统响应延迟。

以上两类调度行为对于调度队列的维护方式要求完全相反，让权行为要求任务尽可能晚地得到下一次调度，而挂起后唤醒行为要求任务尽可能早地得到唤醒后的响应。因此我们设计了双优先级的调度队列，其中 `idle` 调度队列专门用于放置通过让权行为压入的任务，`normal` 调度队列用于放置挂起后唤醒的任务。在调度时会优先选择 `normal` 队列中的任务进行执行。

```

1 pub struct DualPrioScheduler {
2     normal: FifoScheduler,
3     idle: FifoScheduler,
4 }
5 impl Scheduler<Info> for DualPrioScheduler {
6     fn push(&mut self, runnable: Runnable<Info>, info: async_task::ScheduleInfo) {
7         match info.woken_while_running {
8             true => self.idle.push(runnable, info),
9             false => self.normal.push(runnable, info),
10        }
11    }
12    fn pop(&mut self) -> Option<Runnable<Info>> {
13        if let Some(runnable) = self.normal.pop() {
14            return Some(runnable);
15        }
16        self.idle.pop()
17    }
18 }

```

### 2.3.3 分步队列调度

上节中双优先级队列调度方式有一个明显的缺陷：当 `normal` 队列中存在频繁的唤醒行为时，`idle` 队列会因为调度优先级原因一直无法得到执行，从而导致严重的饥饿情况。当 `normal` 队列中的任务依赖于 `idle` 队列中的任务发送的信息时，甚至会导致调度层面上的死锁。

由此，我们借鉴了 Linux 中的 MLFQ 调度算法与 FSCAN 调度算法，设计了较为简易的分步多级调度队列，通过维护 `current` 队列和 `expired` 队列，实现对于优先级队列的轮询式调度。具体而言，调度队列在执行 `push()` 方法的时候，并不会将任务插入到当前的 `current` 队列当中，而是将新的任务插入到 `expired` 队列当中，直到 `current` 队列中所有任务执行完毕，再交换两个队列实现轮询式调度。

该调度设计良好地避免了任务饥饿，同时兼顾了任务的调度优先级。

```
1 type ExpiredSchedulerInnerImpl = DualPrioScheduler;
2 pub struct ExpiredScheduler {
3     current: SyncUnsafeCell<ExpiredSchedulerInnerImpl>,
4     expire: SyncUnsafeCell<ExpiredSchedulerInnerImpl>,
5 }
6 impl ExpiredScheduler {
7     fn switch_expire(&mut self) {
8         core::mem::swap(&mut self.current, &mut self.expire);
9     }
10 }
11 impl Scheduler<Info> for ExpiredScheduler {
12     fn push(&mut self, runnable: Runnable<Info>, info: async_task::ScheduleInfo) {
13         self.expire.as_ref_mut().push(runnable, info);
14     }
15     fn pop(&mut self) -> Option<Runnable<Info>> {
16         let current = self.current.as_ref_mut();
17         let res = current.pop();
18         if let None = res.as_ref() {
19             self.switch_expire();
20             self.current.as_ref_mut().pop()
21         } else {
22             res
23         }
24     }
25 }
```

### 2.3.4 多级调度

部分测例（如：`cyclictest`）对于操作系统调度实时性提出了要求。而上节中的双优先级 + 分步队列调度方式在避免饥饿的同时，也放弃了对于唤醒的高优先级任务的快速响应能力。因此其只适用于常规任务的维护，无法提供良好的实时任务支持。为了实现实时任务的快速响应，我们必须在常规任务队列的基础上再创建额外的实时调度队列。

目前 NoAxiom 只在多级队列中维护了单一优先级的实时调度队列，其内部维护了一个 FIFO 调度队列。当对多级调度队列进行插入操作时，调度器会通过共享裸指针的方式获取当前待调度任务的调度实体 `SchedEntity` 数据以获取其调度优先级 `SchedPrio`，并通过优先级决定最终插入的队列；当对多级队列进行弹出操作时，调度器会优先弹出实时队列当中的任务，从而实现尽可能早的实时任务响应行为。

```

1 type RealTimeSchedulerImpl = FifoScheduler;
2 type NormalSchedulerImpl = ExpiredScheduler;
3 #[repr(align(64))]
4 pub struct MultiLevelScheduler {
5     realtime: RealTimeSchedulerImpl,
6     normal: NormalSchedulerImpl,
7 }
8 impl Scheduler<Info> for MultiLevelScheduler {
9     fn push(&mut self, runnable: Runnable<Info>, info: async_task::ScheduleInfo) {
10         let entity = runnable.metadata().sched_entity();
11         if let Some(entity) = entity {
12             match entity.sched_prio {
13                 SchedPrio::RealTime(_) => self.realtime.push(runnable, info),
14                 _ => self.normal.push(runnable, info),
15             }
16         } else {
17             self.realtime.push(runnable, info);
18         }
19     }
20     fn pop(&mut self) -> Option<Runnable<Info>> {
21         if let Some(runnable) = self.realtime.pop() {
22             return Some(runnable);
23         }
24         self.normal.pop()
25     }
26 }

```

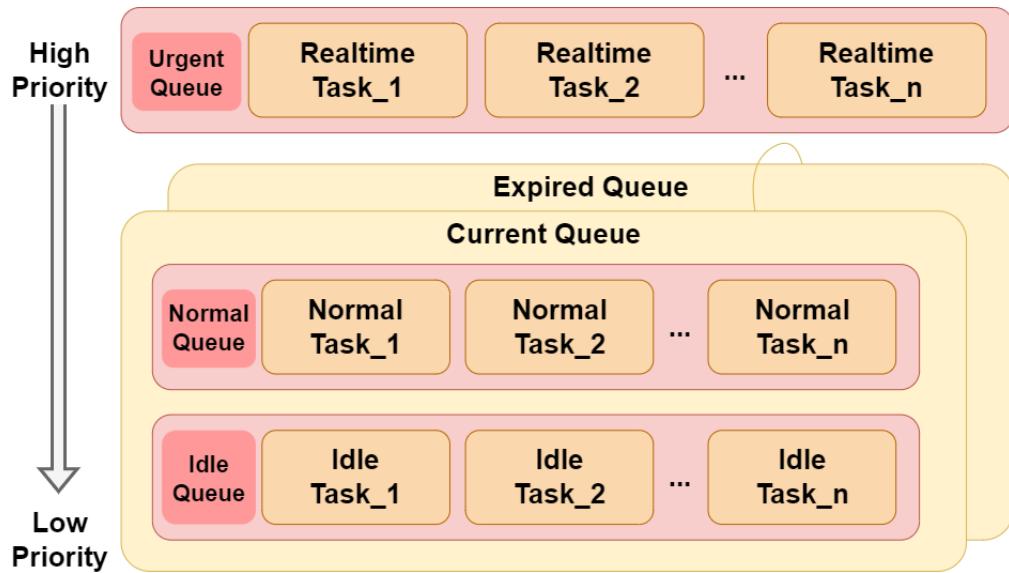


图 4: 多级调度队列设计

## 2.4 实时调度性能调优

在赛事官方提供的测例中，cyclictest 对于系统的实时性提出了要求。该测试程序会创建 200 个高负载的 pipe 读写进程，主动制造惊群效应，从而使得大量被唤醒的 pipe 读写任务堆积在调度器当中等待

调度；同时还会创建时延测试进程，通过调用 `nanosleep` 系统调用测试系统对于挂起后唤醒的实时性进程的响应速度。

NoAxiom 中设计的实时性调度器能够良好地优化实时进程在高负载场景下的响应速度。我们对于 `cyclictest` 进行了实时性调度性能检验，通过比较应用实时调度前后的 `cyclictest` 时延，发现应用调度器的情况下，相较于原有朴素调度器而言均取得了一定的实时性提升。尤其是在高管道负载下的 8 进程调度场景下，响应时延由 48800us 优化至 4750us，调度响应时间提升了 **10 倍**，取得了显著的实时性提升。

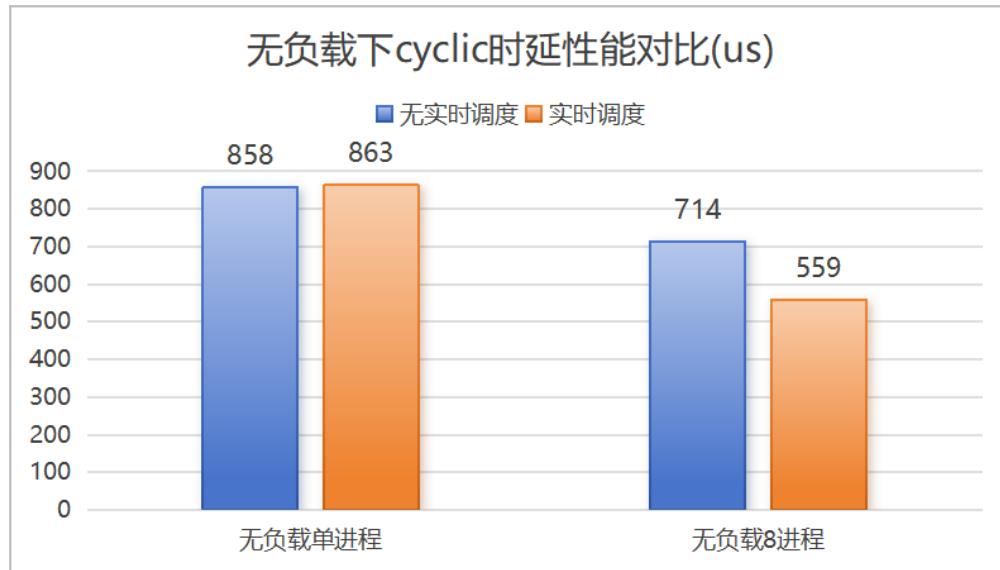


图 5: 无负载下实时性调度性能差异图

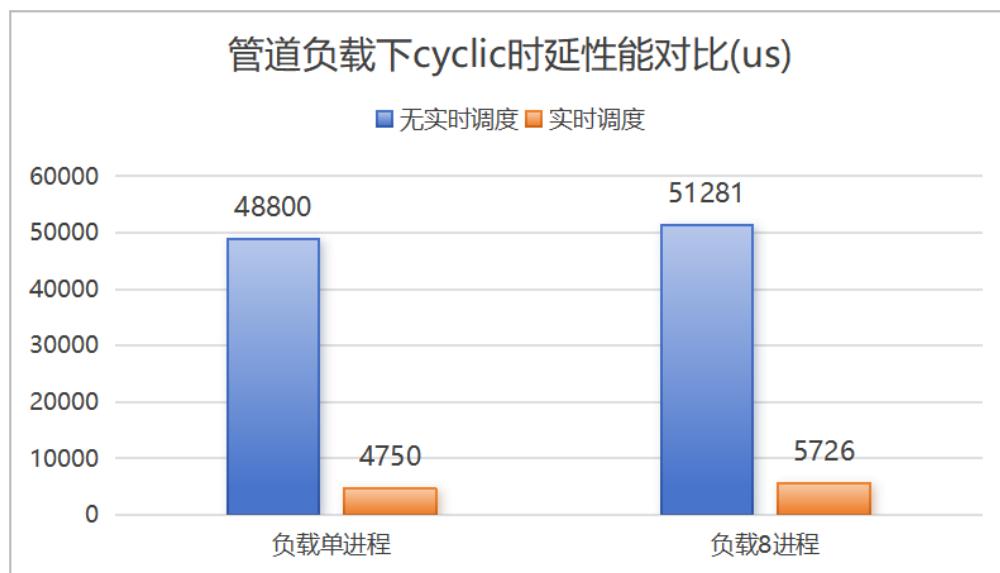


图 6: 管道负载下实时性调度性能差异图

## 2.5 调度信息维护

### 2.5.1 调度实体

为了方便在调度器与任务控制块之间共享调度信息，我们使用调度实体 `SchedEntity` 进行调度相关信息的维护。目前，`SchedEntity` 内维护了任务的调度优先级 `sched_prio`、CPU 掩码 `cpu_mask`、时间统计信息 `time_stat` 等信息，其中就包含了涉及实时性调度的关键成员 `sched_prio`。

```
1 // NoAxiom/kernel/src/sched/sched_entity.rs
2 pub struct SchedEntity {
3     pub sched_prio: SchedPrio, // scheduling priority
4     pub time_stat: TimeInfo,   // task time
5     pub cpu_mask: CpuMask,     // cpu mask
6 }
```

在任务控制块中我们将调度实体设置为了线程独属的数据，并使用 `SyncUnsafeCell` 进行维护。但 Rust 并不允许 `SyncUnsafeCell::get()` 获得的裸指针在调度器与任务控制块之间直接共享，因此我们对调度实体进行了一层 `Sync` 包装，并通过 `SchedMetadata` 进行数据的传递。

```
1 pub struct SchedMetadata {
2     ptr: *mut SchedEntity,
3     tid: usize,
4 }
5 unsafe impl Sync for SchedMetadata {}
6 unsafe impl Send for SchedMetadata {}
```

### 2.5.2 调度优先级

与上节中的描述一致，我们将任务的调度优先级分为三种：

- 实时优先级 `RealTime(usize)`：用于指示特殊的实时任务调度，其中预留了实时任务内部优先级，以应对未来可能的扩展。
- 普通优先级 `Normal`：用于常规执行中因异步操作让权后唤醒的任务，大多数任务均处于这一调度优先级。
- 闲置优先级 `IdlePrio`：我们参照 BFS 和 MuQSS 的设计，为主动让权的进程额外创立了这个调度优先级。由于主动让权的进程通常执行的是轮询行为，其执行通常对时间不敏感，因此将会在执行完普通任务后再进行该优先级任务的执行。

```
1 pub enum SchedPrio {
2     RealTime(usize),
3     Normal,
4     IdlePrio,
5 }
```

## 2.6 异步运行时实现

目前而言，NoAxiom 的异步运行时内部仅维护了一个顶层的多级队列调度器，运行时将通过调用其内部调度器的各类方法，联合外部执行器依赖库创建任务，进一步实现各类运行时功能。

```
1 // NoAxiom/kernel/src/sched/runtime.rs
2 type SchedulerImpl = MultiLevelScheduler;
3 pub struct SimpleRuntime {
4     scheduler: SpinLock<SchedulerImpl>,
5 }
```

### 2.6.1 异步任务的生成

我们通过 `Runtime::spawn()` 方法将一个 `Future` 异步函数对象与任务控制块 `Arc<Task>` 打包为执行器库中的 `Runnable` 对象进行调度。具体而言，当生成一个新任务时，将通过调用 `async_task` 库中的 `spawn()` 方法以获取生成器生成的任务对象 (`Runnable, handle`)，并调用 `Runnable.schedule()` 方法将其投入调度队列中。

我们使用的外部执行器依赖库 `async_task` 中默认关闭了对于其调度对象 `Runnable<T = ()>` 的额外元信息维护。想要通过 `SchedMetadata` 在执行器环境中共享调度信息，应当在创建异步任务时额外调用 `async_task::Builder::new().metadata()` 中的构建方法以携带额外的元信息。此后，即可在调度器中通过调用 `Runnable::metadata()` 方法获取到构建的元信息了。

```
1 // NoAxiom/kernel/src/sched/runtime.rs
2 impl Runtime<SchedulerImpl, Info> for SimpleRuntime {
3     fn spawn<F>(self: &'static Self, future: F, task: Option<&Arc<Task>>)
4         where
5             F: core::future::Future<Output: Send + 'static> + Send + 'static,
6         {
7             let metadata = task
8                 .map(|task| SchedMetadata::from_task(task))
9                 .unwrap_or_else(SchedMetadata::default);
10            let (runnable, handle) = Builder::new().metadata(metadata).spawn(
11                move |_ | future,
12                WithInfo(move |runnable, info| self.schedule(runnable, info)),
13            );
14            runnable.schedule();
15            handle.detach();
16        }
17 }
```

### 2.6.2 异步任务的调度与运行

我们通过 `Runtime::schedule()` 进行单个任务的调度申请，该方法通过调用 `Scheduler::push()` 方法实现向调度器中添加一个新的可执行任务的功能。

我们通过 `Runtime::run()` 进行单个任务的单次执行，该函数会通过调用 `Scheduler::pop()` 函数获取一个待执行的任务，并通过调用异步执行器的 `Runnable::run()` 方法在同步控制流中实现单次异步函数的轮询行为。

```
1 impl Runtime<SchedulerImpl, Info> for SimpleRuntime {
2     fn run(&self) {
3         let runnable = self.scheduler.lock().pop();
4         if let Some(runnable) = runnable {
5             runnable.run();
6         }
7     }
8     fn schedule(&self, runnable: Runnable<Info>, info: async_task::ScheduleInfo) {
9         self.scheduler.lock().push(runnable, info);
10    }
11 }
```

我们的内核一旦初始化完成，就将持续访问运行时执行任务。该行为被包装为 `run_task()` 函数，作为最顶层的执行循环，实现运行时的不间断执行行为。

```
1 #[no_mangle]
2 pub fn run_tasks() -> ! {
3     loop {
4         timer_handler();
5         RUNTIME.run();
6     }
7 }
```

## 三、进程管理

### 3.1 任务控制块

在操作系统中，进程是资源分配的单位，其拥有独立的资源空间。而线程是共享进程资源的执行单元，作为属于某一进程下的活动对象参与执行。在 Linux 系统中，进程和线程被统一表示为轻量级进程（Lightweight Process, LWP），并通过统一的任务控制块进行进程与线程资源的管理。在 NoAxiom 当中，我们延续了这一设计，将任务抽象为 `Task` 对象，从而实现了对于进程和线程的统一管理。

#### 3.1.1 资源访问属性

在任务控制块中，资源的访问属性尤为重要，错误的访问级别会导致锁和裸指针的错误使用，进而导致频繁的资源竞争或多核间数据一致性问题。为了清晰表示任务资源的访问属性，在设计任务控制块的过程当中，我们将任务内部资源按照是否共享、是否可访问、是否可变大致分为下面四种类型。

```
1 // shared between threads
2 type SharedMut<T> = Arc<SpinLock<T>>;
3
4 // mutable resources mostly used in current thread
5 // it could be accessed by other threads through process manager
6 // so protect it with spinlock
7 type Mutable<T> = SpinLock<T>;
8
9 // only used in current thread, mutable resources without lock
10 // SAFETY: these resources won't be shared with other threads
11 type ThreadOnly<T> = SyncUnsafeCell<T>;
12
13 // read-only resources, could be shared safely through threads
14 type Immutable<T> = T;
```

上述定义的具体解释如表 2 所示。

表 2: 进程内部成员访问属性定义

字段名称	访问属性	具体含义	具体维护策略
<code>SharedMut&lt;T&gt;</code>	共享、可变	线程间共享资源	使用引用计数和自旋锁维护
<code>Mutable&lt;T&gt;</code>	可访问、可变	可被管理器访问的独占资源	使用自旋锁维护
<code>ThreadOnly&lt;T&gt;</code>	独享、可变	仅当前线程可访问的独占资源	使用裸指针维护
<code>Immutable&lt;T&gt;</code>	不可变	只读数据，可被任意读取	直接保存在控制块中

#### 3.1.2 任务控制块

基于以上访问权限设计，我们可以给出任务控制块的具体内部成员定义，用于记录任务相关的所有信息，如下所示。其具体字段的含义在表 6 中给出。

```
1 // NoAxiom/kernel/src/task/task.rs
```

```

2 #[repr(C, align(64))]
3 pub struct Task {
4     // mutable
5     pcb: Mutable<PCB>, // task control block inner, protected by lock
6
7     // thread only / once initialization
8     tcb: ThreadOnly<TCB>, // thread control block
9     cx: ThreadOnly<TaskContext>, // trap context
10    sched_entity: ThreadOnly<SchedEntity>, // sched entity for the task
11    waker: Once<Waker>, // waker for the task
12    ucx: ThreadOnly<UserPtr<UContext>>, // ucontext for the task
13
14     // immutable
15     tid: Immutable<TidTracer>, // task id, with lifetime holded
16     tgid: Immutable<TGID>, // task group id, aka pid
17     tg_leader: Immutable<Once<Weak<Task>>>, // thread group leader
18
19     // shared
20     fd_table: SharedMut<FdTable>, // file descriptor table
21     cwd: SharedMut<Path>, // current work directory
22     sa_list: SharedMut<SigActionList>, // signal action list, saves signal handler
23     memory_set: SharedMut<MemorySet>, // memory set for the task
24     thread_group: SharedMut<ThreadGroup>, // thread group
25     pgid: Arc<AtomicUsize>, // process group id
26     futex: SharedMut<FutexQueue>, // futex wait queue
27     itimer: SharedMut<ITimerManager>, // interval timer
28 }

```

有关任务控制块内部关键成员的说明如下。

1. 任务 ID：我们使用全局的 TID 分配器进行任务 ID 的原子性分配，并基于 RAI<sup>I</sup> 的思想，通过 `TidTracer` 将其生命周期与任务绑定。TID 将作为任务的唯一标识参与各类管理器的访问中。
2. 进程控制块：为了保证原子性地维护进程间关系以及信号信息，我们选择将所有进程间父子关系和信号控制信息全部放置在进程控制块 PCB 当中进行维护，并使用自旋锁对其访问原子性进行保护。下面是 PCB 的成员，其具体字段含义在表 4 中给出。
3. 共享资源：以 `clone(2)`<sup>5</sup> 为代表的若干系统调用会通过设置 `CloneFlags` 的方式，规定进程在复制时的资源共享情况。这些资源包括：`fd_table`, `cwd`, `sa_list`, `memory_set`, `thread_group`, `pgid`, `futex`, `itimer` 等。通过这些共享资源，当前线程可以与其他兄弟进程或线程共享资源，从而实现更高效的资源利用。为了维护共享资源，我们使用引用计数指针 `Arc<T>` 对其进行生命周期维护，并使用自旋锁 `SpinLock` 维护修改的原子性。`SpinLock`<sup>6</sup> 是一种自旋锁，通过我们自定义的 `NoIrqLockAction` 机制，实现持有锁期间的中断使能关闭机制，因而它适用于实现短时锁定时的原子性保证，尤其是内核中对延迟敏感的部分。而共享引用计数 `Arc` 允许线程之间所有权的共享，在引用计数变为 0 时数据才会被释放，实现了共享资源生命周期的维护。
4. 独占资源：我们保证当前任务独占资源不会被其他任务以任何形式单独访问，因此这些资源不需要额外的锁开销进行一致性维护，可以使用裸指针直接进行访问。由于裸指针包装

<sup>5</sup><https://www.man7.org/linux/man-pages/man2/clone.2.html>

<sup>6</sup><https://github.com/os-module/kernel-sync>

器 `UnsafeCell` 未实现 `Sync` 属性，无法安全地在进程间传递，因此我们使用实现了 `Sync` 属性的 `SyncUnsafeCell` 绕过编译器检查，允许其在多线程环境下安全地传递。

- 资源访问：通过 `Arc<SpinLock>` 和 `SyncUnsafeCell` 的严格定义，任务控制块保护了共享数据访问原子性以实现线程安全，同时也保证了独占资源高效访问的高并发性，二者结合实现了并发情况下的高效性与正确性。

表 3: 任务控制块成员定义

字段名称	访问属性	具体含义
<code>pcb</code>	可变	进程控制块，用于维护进程信息原子性修改
<code>tcb</code>	独占	线程控制块，用于维护其他线程独占信息
<code>cx</code>	独占	任务上下文，保存任务的上下文信息
<code>sched_entity</code>	独占	调度实体，包含任务的调度信息
<code>waker</code>	独占	唤醒器，用于异步任务唤醒
<code>ucx</code>	独占	信号处理时保存的用户态上下文信息
<code>tid</code>	不可变	任务 ID，同时绑定了 tid 的生命周期
<code>tgid</code>	不可变	线程组 ID，也用作表示进程 ID 号
<code>tg_leader</code>	不可变	保存对线程组先导任务的弱引用
<code>fd_table</code>	共享可变	文件描述符表
<code>cwd</code>	共享可变	当前工作目录
<code>sa_list</code>	共享可变	信号处理器列表
<code>memory_set</code>	共享可变	内存空间，保存任务的全部内存映射信息
<code>thread_group</code>	共享可变	线程组，维护当前线程组中的任务引用
<code>pgid</code>	共享可变	进程组 ID，使用原子操作维护可变性
<code>futex</code>	共享可变	Futex 等待队列
<code>itimer</code>	共享可变	间隔定时器管理器

表 4: 进程控制块成员定义

字段名称	具体含义
<code>status</code>	任务状态，表示任务的当前状态
<code>exit_code</code>	任务退出码，表示任务的退出状态
<code>children</code>	子任务列表，保存当前任务的子进程的强引用
<code>parent</code>	父任务引用，保存当前任务的父进程的弱引用
<code>pending_sigs</code>	待处理信号，保存当前任务的待处理信号信息
<code>sig_stack</code>	信号备用栈，保存信号处理时的备用栈信息
<code>robust_list</code>	保存 robust 列表信息

## 3.2 任务状态

### 3.2.1 任务状态转化关系

我们在设计任务状态时，参考 Linux 的任务状态设计了以下若干任务状态，并将任务状态保存在 `PCB::TaskStatus` 当中，以自旋锁进行互斥访问。

任务状态转换如图 7 所示。

1. `Running`：当前任务正在处于运行状态，正在被执行器执行中。
2. `Runnable`：当前任务可运行但未被执行，存在于调度器中。
3. `Suspend`：当前任务挂起，保存在反应器中等待外部事件唤醒。
4. `Stopped`：当前任务被信号标记为暂停，检测到该状态后将挂起（`Suspend`）当前任务。
5. `Terminated`：当前任务被 `exit` 等系统调用标记为退出。当任务检测到处于该状态后会执行任务退出处理例程并进入僵尸态（`Zombie`）。
6. `Zombie`：当前任务执行完毕并进入僵尸态，等待父进程对其资源进行回收。

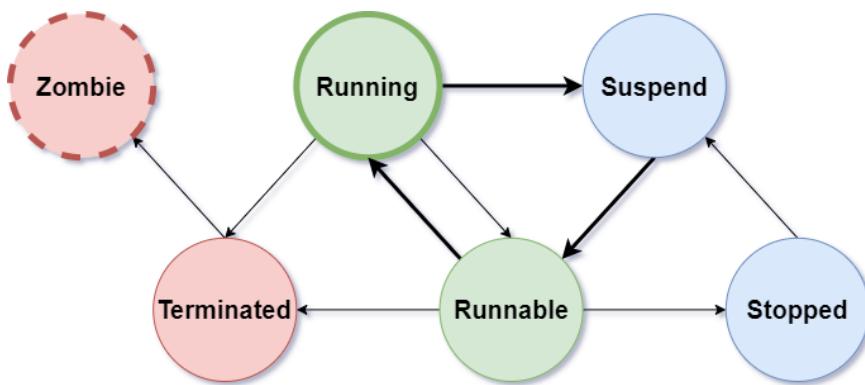


图 7: 任务状态转换状态机

### 3.2.2 调度相关状态简化设计

在查看外部库 `async-task` 的过程当中，我们注意到：该外部库已经通过原子标志位维护了任务的运行（`Running`）、等待（`Runnable`）、挂起（`Suspend`）三个状态了。具体而言，该外部库会通过 `Header` 中的 `state` 成员进行原子访存，以此维护当前任务在调度方面的状态更新原子性。

```
1 // async-task/src/header.rs
2 pub(crate) struct Header<M> {
3     /// Current state of the task.
4     /// Contains flags representing the current state and the reference count.
5     pub(crate) state: AtomicUsize,
6 }
```

例如，在执行 `Runnable.run()` 将某个任务从等待态变为运行态的时候，该外部库会通过原子访存将 `SCHEDULED` 位重置，以通知运行时该任务可再次通过 `Runnable.schedule()` 加入调度队列中。

```
1 let state = (*raw.header).state.fetch_and(!SCHEDULED, Ordering::AcqRel);
```

因此在 NoAxiom 中，我们将常规的三种执行状态 `Running`、`Runnable`、`Suspend` 实际统一为固定的 `Normal` 状态，而实际的执行状态则通过任务在运行时中保存的调度模块决定。有关具体保存的调度模块，详见 2.1.2 章节中的描述，其中图 2 详细描述了本章节中的对应状态转化图。

## 3.3 任务间关系

### 3.3.1 进程间关系

在操作系统中，进程通过 `clone(2)` 等系统调用创建子进程，以此构成进程之间的父子关系。当子进程退出时，为了保证进程资源的一致性，其占用的资源需要父进程帮助回收。具体而言，任务控制块的 `pcb` 中维护了进程的父子关系信息，通过持有子进程的强引用 `children: Vec<Arc<Task>>` 防止其提早析构，并持有父进程的弱引用 `parent: Option<Weak<Task>>` 便于临时获取父进程的任务控制块。

### 3.3.2 线程间关系

我们参照 Linux 中的设计，使用线程组（Thread Group）维护同一进程下的线程。进程组中将维护一个主线程，作为当前线程组的进程级任务，负责管理共享资源以及所有同一线程组下的线程级任务；而其余的所有任务均被视作线程，线程之间不存在父子关系，也不需要对进程资源的维护直接负责。当进程退出时，其他所有的线程都需要退出，并等待父进程回收其资源；当非主线程退出时，将会主动进行资源释放。线程组结构如图 8 所示。



图 8: 线程组结构

具体而言，每个进程下的所有线程都会通过 `Task::thread_group` 共享同一个线程组。线程组内部维护了从任务号 `tid` 到任务控制块 `Weak<Task>` 的映射，用于线程组信息的快速访问与维护。此外，任务控制块内会额外维护一个主线程 `tg_leader` 用于代表整个进程，该主线程的任务号 `tid` 也将作为线程组号/进程号 `tgid` 保存在任务控制块中。通过访问线程组，可以较为便利地对同一进程下的所有线程信息进行读取和维护。

```
1 pub struct ThreadGroup(pub BTreemap<TID, Weak<Task>>);
2 impl ThreadGroup {
3     pub const fn new() -> Self {
4         Self(BTreemap::new())
5     }
6     pub fn insert(&mut self, task: &Arc<Task>) {
```

```

7     self.0.insert(task.tid(), Arc::downgrade(&task));
8 }
9 pub fn remove(&mut self, taskid: usize) {
10    self.0.remove(&taskid);
11 }
12 }
```

### 3.3.3 任务管理器

仅仅通过进程父子关系、线程组关系的维护并不足以实现完整的任务管理。我们参照 Linux 的设计，创建了任务管理器 TaskManager 用于实现全局的任务关系维护。与线程组类似，任务管理器将保存任务标识 tid 到任务控制块 Weak<Task> 的映射关系，以用于所有任务信息的快速访问与维护。

```

1 /// NoAxiom/kernel/src/task/manager.rs
2 pub struct TaskManager(pub SpinLock<BTreesMap<TID, Weak<Task>>>);
3 pub static TASK_MANAGER: Lazy<TaskManager> = Lazy::new(|| TaskManager::new());
```

### 3.3.4 进程组管理器

在兼容于 POSIX 标准的操作系统中，进程组（Process group）是指一个或多个进程的集合。进程组被使用于控制信号的分配。对于一个进程组发出的的信号，会被个别递送到这个组群下的每个进程成员中。<sup>7</sup>NoAxiom 严格遵循 POSIX 相关标准，对于进程组使用进程组管理器 ProcessGroupManager 进行维护。其内部通过维护同一进程组号 PGID 下的所有进程，实现某一进程组信息的快速访问与修改。

```

1 pub struct ProcessTracer {
2     pid: PID,
3     proc: Weak<Task>,
4 }
5 pub struct ProcessGroupManager(BTreesMap<PGID, Vec<ProcessTracer>>);
6 pub static PROCESS_GROUP_MANAGER: Lazy<SpinLock<ProcessGroupManager>> =
7     Lazy::new(|| SpinLock::new(ProcessGroupManager::new()));
```

## 3.4 任务生命周期

任务的生命周期由三个主要步骤组成：任务的创建，任务的执行，以及任务的退出。我们在图 9当中给出了详细的任务生命周期及其运行流程的图示。

---

<sup>7</sup>[https://en.wikipedia.org/wiki/Process\\_group](https://en.wikipedia.org/wiki/Process_group)

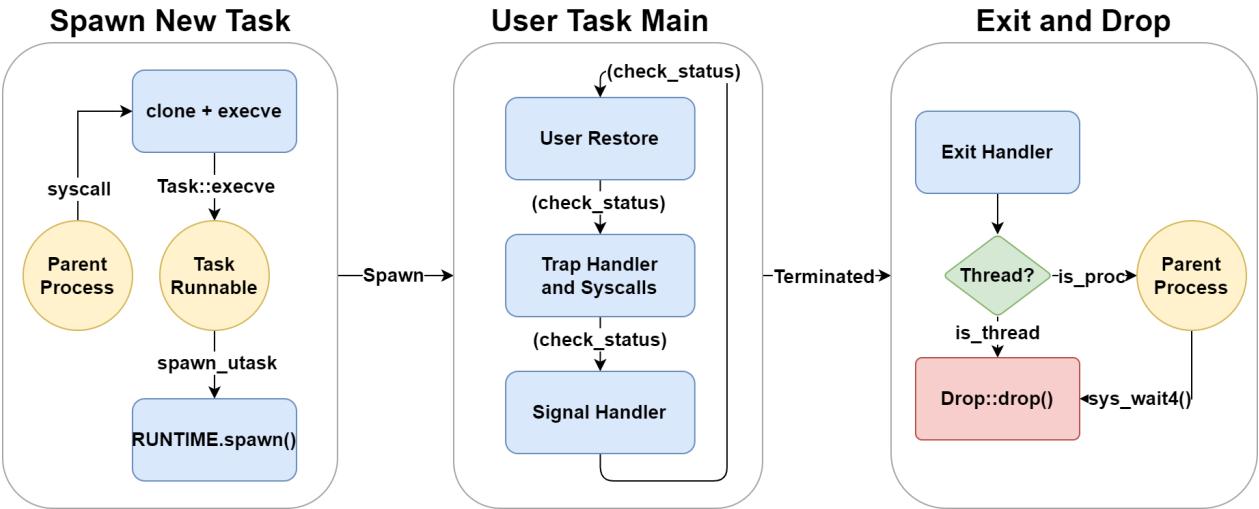


图 9: 任务生命周期

## 3.5 任务的创建

### 3.5.1 初始进程的加载

在 NoAxiom 当中，初始进程以内嵌的二进制形式 .incbin 随内核一同加载到内存当中。在编译内核时，我们会提前编译好初始进程所需的 elf 文件，并通过 build.rs 构建到内核代码的 link\_apps.S。

例如，当选择的初始进程名为 run\_busybox 时，我们生成的部分汇编代码如下：

```

1 .section .data
2 .global run_busybox_start
3 .global run_busybox_end
4 .align 3
5 run_busybox_start:
6     .incbin "../NoAxiom-OS-User/bin/run_busybox"
7 run_busybox_end:

```

在内核当中，我们通过自定义宏 use\_apps! 通过编译选项自动选择初始进程名，并将对应的程序段链接到内核的可执行文件内。通过这种方式，NoAxiom 可以快速修改用户程序的功能实现，并快速切换到所需的初始进程。

```

1 macro_rules! use_app {
2     ($name:literal) => {
3         extern "C" {
4             #[link_name = concat!($name, "_start")]
5             fn app_start();
6             #[link_name = concat!($name, "_end")]
7             fn app_end();
8         }
9         pub const INIT_PROC_NAME: &str = $name;
10    };
11 }
12 #[cfg(feature = "busybox")]
13 use_app!("run_busybox");
14 #[cfg(feature = "runitests")]

```

```
15 use_app!("run_tests");
```

由于初始进程的数据一开始位于内核数据段而不是文件系统当中，为了方便文件系统对于块设备的统一管理，复用从文件加载 elf 的代码，我们选择先通过 `File::write_at` 方法将内嵌的数据写入文件系统中，再通过创建的初始进程路径进行数据加载。

```
1 pub fn schedule_spawn_with_path() {
2     spawn_ktask(async move {
3         let path_str = format!("/{}", INIT_PROC_NAME);
4         let path = Path::from_or_create(path_str, InodeMode::FILE).await;
5         let file = path.dentry().open().unwrap();
6         file.write_at(0, get_file()).await.unwrap();
7         let elf_file = path.dentry().open().unwrap();
8         let elf = MemorySet::load_elf(&elf_file).await.unwrap();
9         let task = Task::new_process(elf).await;
10        spawn_utask(task);
11    });
12 }
```

### 3.5.2 普通任务的加载

除了初始进程外的所有其他进程均依赖于 `clone(2)`、`execve(2)` 等系统调用来创建新任务。当执行 `execve(2)` 时，任务应当将当前执行的可执行文件镜像替换为另一文件镜像，初始化用户栈信息，更新用户上下文，关闭相应的任务资源，并终止其他线程运行。

```
1 pub async fn execve(
2     self: &Arc<Self>,
3     path: Path,
4     args: Vec<String>,
5     envs: Vec<String>,
6 ) -> SysResult<()> {
7     let elf_file = path.dentry().open()?;
8     let ElfMemoryInfo {
9         memory_set,
10        entry_point,
11        user_sp,
12        mut auxs,
13    } = MemorySet::load_elf(&elf_file).await?;
14     memory_set.memory_activate();
15     self.terminate_threads();
16     self.change_memory_set(memory_set);
17     let (user_sp, _argc, _argv_base, _envp_base) =
18         self.init_user_stack(user_sp, args, envs, &mut auxs);
19     *self.trap_context_mut() = TrapContext::app_init(cx(entry_point, user_sp));
20     self.sa_list().reset();
21     self.fd_table().close_on_exec();
22     Ok(())
23 }
```

## 3.6 任务运行流程

任务的运行是任务生命周期中占比最大的部分，任务将在用户态和内核态之间持续切换，并根据用户态异常的种类，提供系统调用服务或其他异常处理功能。任务的主要执行流程如图 10 所示。我们通过调用 `UserTaskFuture` 内部包装的异步函数来实现用户程序的执行。

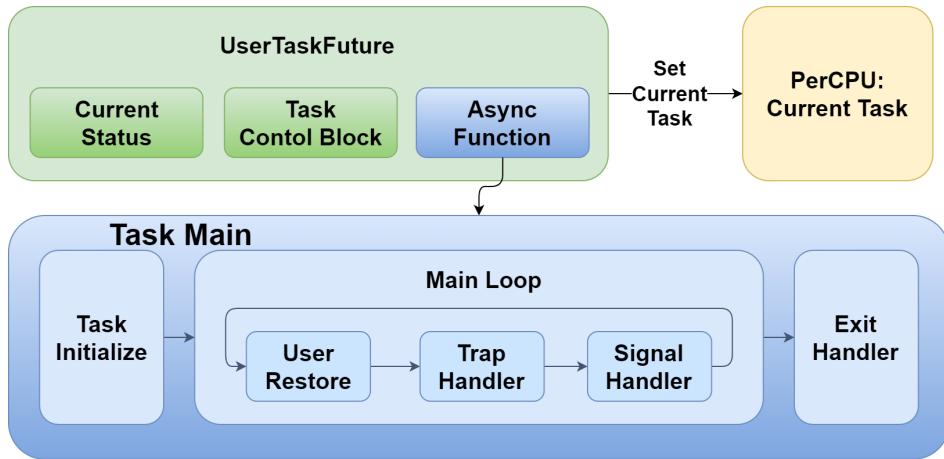


图 10: 任务执行流程

### 3.6.1 任务切换

NoAxiom 的所有任务均运行在异步运行时下，所有用户任务均被包装在异步的 `UserTaskFuture` 当中，等待执行器执行。在用户进程的生命周期内，其所有执行功能均在 `UserTaskFuture` 内部完成。

`UserTaskFuture` 负责任务切换流程，当切入任务时，`UserTaskFuture` 将会执行以下步骤：

1. 关闭中断，保证任务切入的原子性。
2. 记录当前任务的切入时间，便于进行时间统计和时间片维护。
3. 访问当前核心的 `per_cpu` 字段，设置当前运行的任务控制块。
4. 激活任务的地址空间，刷新 TLB 内缓存的地址翻译信息。
5. 恢复任务在切出前记录的中断使能信息。

而在切出任务时，`UserTaskFuture` 也将镜像地执行上述行为，从而实现了任务信息切入切出的原子性与完整性。

```
1 pub struct UserTaskFuture<F: Future + Send + 'static> {
2     pub task: Arc<Task>,
3     pub future: F,
4 }
5 impl<F: Future + Send + 'static> Future for UserTaskFuture<F> {
6     type Output = F::Output;
7     fn poll(self: Pin<&mut Self>, cx: &mut Context<'_>) -> Poll<Self::Output> {
8         Arch::disable_interrupt();
```

```

9     let this = unsafe { self.get_unchecked_mut() };
10    let task = &this.task;
11    let future = &mut this.future;
12    task.time_stat_mut().record_switch_in();
13    current_cpu().set_task(task);
14    fence(Ordering::AcqRel);
15    task.restore(cx_int_en());
16    // ===== before executing task future =====
17    let ret = unsafe { Pin::new_unchecked(future).poll(cx) };
18    // ===== after executing task future =====
19    task.record(cx_int_en());
20    Arch::disable_interrupt();
21    task.time_stat_mut().record_switch_out();
22    current_cpu().clear_task();
23    fence(Ordering::AcqRel);
24    Arch::enable_interrupt();
25    ret
26 }
27 }
```

### 3.6.2 任务主函数

上节中的 `UserTaskFuture` 只负责任务信息的切入与切出，而实际的用户程序执行则在其内部包装的异步函数 `task_main` 中发生。异步函数 `task_main()` 负责任务主函数运行、用户上下文切换、任务退出等关键流程的执行，是真正意义上用于用户程序运行的主函数。

为了便于通过控制块进行任务唤醒，在进入用户程序的执行循环前，`task_main` 将通过 `take_waker` 获得当前异步函数唤醒器并注册在任务控制块中。

在初始化完任务控制信息后，`task_main` 将进入用户程序主循环，用于持续的内核/用户态切换，并向用户提供系统调用等服务功能。具体而言，用户程序主循环的执行流程如下：

1. 恢复用户态执行：通过调用硬件抽象层的 `trap_restore()` 函数，恢复用户的执行上下文。此外，在进入用户态和返回内核态的过程中，将统计用户态的执行时间信息。
2. 处理用户态例外：通过调用硬件抽象层的 `read_trap_type`，获取用户态发生的例外类型，并进入例外处理例程进行处理。该处理例程将提供系统调用、懒分配、写时复制、外部中断处理、时间片让权等关键内核功能的服务。
3. 信号处理：检查是否存在挂起的未处理信号，若存在，则执行信号处理例程。

此外，在每个执行流程之间，内核都将对当前任务状态进行检查，如果检测到当前任务已经被标记为退出（`TaskStatus::Terminated`），则将立即退出用户程序的循环执行，并执行退出处理例程 `exit_handler()`。如果检测到当前任务被信号标记为暂停（`TaskStatus::Stopped`），则将立即挂起当前进程。

```

1 pub async fn task_main(task: Arc<Task>) {
2     task.set_waker(take_waker().await);
3     loop {
4         // kernel -> user, restore context and return to user mode
5         task.time_stat_mut().record_trap_in();
```

```

6     let cx = task.trap_context_mut();
7     Arch::trap_restore(cx);
8     let trap_type = Arch::read_trap_type(Some(cx));
9     task.time_stat_mut().record_trap_out();
10    match task.pcb().status() {
11        TaskStatus::Terminated => break,
12        TaskStatus::Stopped => suspend_now().await,
13        _ => {}
14    }
15    // user -> kernel, enter the handler
16    user_trap_handler(&task, trap_type).await;
17    match task.pcb().status() {
18        TaskStatus::Terminated => break,
19        TaskStatus::Stopped => suspend_now().await,
20        _ => {}
21    }
22    // check signal before return to user
23    task.check_signal(None).await;
24    match task.pcb().status() {
25        TaskStatus::Terminated => break,
26        TaskStatus::Stopped => suspend_now().await,
27        _ => {}
28    }
29 }
30 task.exit_handler().await;
31 }
```

### 3.7 任务资源释放

当任务在 `task_main` 函数中检测到当前状态为 `Terminated` 之后，将跳出用户程序的执行过程并执行退出处理例程 `exit_handler`。这一过程标志着任务生命周期进入尾声。

Rust 语言对于对象的生命周期做出了严格的规定，当对象生命周期结束时，将自动调用其析构方法。我们使用 RAII 的思想，使用 Rust 原生的 `Arc<T>` 原子引用计数指针包装任务控制块对象 `Task`，这是一个线程安全的共享指针，其内部将使用原子计数器 `AtomicUsize` 维护当前指针被引用的次数。当 `Arc` 指针的引用计数归零时，表明其生命周期结束，将会自动调用 `Task` 的 `Drop::drop()` 函数以释放内部资源。

对于线程而言，其资源应当由自身释放。当执行完毕线程的退出处理例程之后，线程的任务控制块 `Arc<Task>` 的引用计数归零，自动执行其 `Drop` 方法释放内部所有的线程独占资源。

对于进程而言，为了防止资源提前释放导致的资源错误复用问题，其内部资源应当延迟回收，由父进程负责通过 `wait4()` 等系统调用进行子进程资源的回收。我们已经在程序控制块 `pcb` 内部保存了每个进程的子进程的强引用 `children: Vec<Arc<Task>>`，可以避免子进程因为引用计数归零导致的提前释放。在检测到子进程退出后，父进程会将子进程的强引用从 `children` 字段中删除，使其引用计数清零而自动触发任务控制块的 `Drop` 方法。

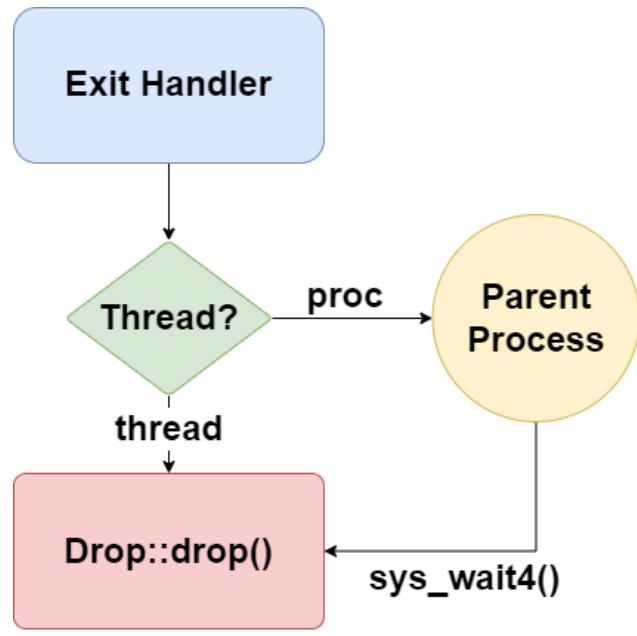


图 11: 任务资源释放流程

## 四、内存管理

### 4.1 地址空间结构

#### 4.1.1 内核地址空间初始化

目前主流的地址空间设计大致分为共享地址空间和隔离地址空间两类。隔离的地址空间设计拥有更高的安全性，能够一定程度上防止侧信道攻击等针对缓存的攻击方式，但这种设计方式在切换用户/内核上下文时，需要频繁切换根页表并冲刷 TLB，导致上下文切换的性能下降。

考虑到系统性能需求，NoAxiom 选择将用户与内核共享地址空间，这种维护方式无需频繁切换根页表，有效避免了 TLB 冲刷开销。同时，为了在同一地址空间下隔离用户与内核空间，我们参照了 Linux 中的设计，选择将用户地址设置在低半地址空间，内核地址设置在高半地址空间。

具体而言，我们在硬件抽象层中根据具体架构定义了 `KERNEL_ADDR_OFFSET` 用于指示具体的高半地址空间前缀。这一前缀在 RISC-V64 架构中为 `0xffff_ffc0_0000_0000`，而这前缀在 LoongArch64 架构中为 `0x9000_0000_0000_0000`。

对于 RISC-V64 架构而言，RV64 硬件级页表翻译使得我们能够便利地创建大页，因此 NoAxiom 将在系统启动时使用大页维护的临时页表进行启动维护，其内部将整个内核区域不区分读写权限地映射到了高半地址空间。在随后初始化过程中，再创建内核对应的地址空间 `KERNEL_SPACE: MemorySet`，将具体的内核程序段映射到 `0xffff_ffc0` 开头的高半地址空间后，再激活真正的内核地址空间。

RV64 中，在系统一开始启动时将会把页表设置给定的临时页表，并启动页表翻译功能。

```
1 la      t0, {page_table}
2 srlt    t0, t0, 12
3 li      t1, 8 << 60
4 or      t0, t0, t1
5 csrw    satp, t0
6 sfence .vma
```

对于 LoongArch64 架构而言，由于 LA64 架构支持了地址映射窗口，NoAxiom 选择直接配置 CSR 寄存器中 `DMWIN1` 寄存器，将其地址映射窗口置为 `0x9000`。由于在 LoongArch 架构中，在 TLB 进行地址翻译时，将会优先使用命中地址映射窗口的物理地址，因此也无需额外进行内核空间的页表初始化。

```
1 ori      $t0, $zero, 0x11      # CSR_DMWIN1_MAT | CSR_DMWIN1_PLVO
2 lu52i.d   $t0, $t0, -1792     # CA, PLVO, 0x9000 xxxx xxxx xxxx
3 csrwr    $t0, 0x181         # LOONGARCH_CSR_DMWIN1
```

关于高半地址空间在系统启动时的具体维护事项，详见 5.3 中的描述。

#### 4.1.2 内核地址空间排布

内核地址空间的排布较为简单，其地址空间按照内核可执行文件给定的程序段，从起始地址开始依次排布。而剩余的地址空间则用于进行物理页帧的分配。

- 我们给定 `ekernel` 用于表示内核 elf 的结束地址。从起始地址到 `ekernel` 的数据均位于内核 elf 中，我们将通过内核 elf 给定的读写属性进行数据映射。

- 内核的 `bss` 段中将有部分空间用于维护内核堆空间，目前 QEMU 中的内核堆空间大小为 48MB。
- 从 `ekernel` 开始一直到物理地址结束的空间将全部用于物理页帧分配，在现在的 NoAxiom 中，大约能够支持从内存中分配 19000 到 20000 个物理页帧。
- 为了避免用户更改内核内容，内核空间完全不支持用户态访问，我们选择将内核中的跳板页映射在用户空间中。

#### 4.1.3 用户地址空间分布

由于启用了页表翻译功能，我们的用户地址空间可以较为宽松地进行配置，为不同类型的用户空间分配充足的虚地址空间。具体而言，以 RISC-V64 架构为例，我们的用户空间映射策略如表 5 所示，其中按照地址从低到高以此给出了所有地址映射基址，映射长度及其对应的功能。

表 5: 用户空间映射区域表

功能	映射基址	(最大) 映射长度	默认映射属性
用户程序 elf	由用户程序决定	由用户程序决定	由用户程序决定
用户栈 stack	<code>ELF_END + 0x1000</code>	<code>0x80_0000</code>	R, W
用户堆 brk	<code>STACK_END + 0x1000</code>	<code>0x753_0000</code>	R, W
共享内存 shm	<code>0x5_0000_0000</code>	由系统调用决定	由系统调用给定
内存映射 mmap	<code>0x6_0000_0000</code>	<code>0x1000_0000</code>	由系统调用给定
信号跳板页	<code>0x1F_FFFF_F000</code>	<code>0x1000</code>	R, X

这里需要特别注意，与部分往届基于 RISC-V 架构的操作系统不同，NoAxiom 选择将 mmap 区域的基址设置为更高的 `0x6_0000_0000`，以防止在 LoongArch 架构下与用户 elf 程序段起始地址 `0x1_2000_0000` 发生重叠而导致地址映射错误。

此外，我们专门在位于 `USER_MEMORY_END - PAGE_SIZE`，即用户映射区域尾段的位置映射了信号跳板页，专门用于执行内核提供的信号处理相关系统调用及其相关上下文切换操作，如 `sigreturn` 系统调用等，该页面由内核固定映射，每次加载内存区域时都会进行映射和相关指令的加载。

## 4.2 内存管理器

### 4.2.1 堆分配管理器

`no_std` 环境下的 Rust 要求应当给定一个全局堆分配器用于进行堆空间申请。我们使用开源外部依赖库 `buddy_system_allocator` 作为 NoAxiom 的全局堆分配器，该依赖库使用伙伴算法对堆内数据进行维护。同时为了保证堆分配的原子性，我们使用关中断的 `SpinLock` 保护堆分配的过程。

此外，在 NoAxiom 当中，我们在 `.bss` 程序段预留了大小为 48MB 的堆空间用于堆分配器的初始化。当系统初始化的时候，将会对于堆空间调用 `init` 方法进行初始化。

```

1 #[global_allocator]
2 static HEAP_ALLOCATOR: HeapAllocator = HeapAllocator::empty();
3 struct HeapAllocator(SpinLock<Heap<32>>);
```

```

4 unsafe impl GlobalAlloc for HeapAllocator {
5     unsafe fn alloc(&self, layout: Layout) -> *mut u8 {
6         self.0.lock().alloc(layout)
7             .ok().map_or(0 as *mut u8, |allocation| allocation.as_ptr())
8     }
9     unsafe fn dealloc(&self, ptr: *mut u8, layout: Layout) {
10        self.0.lock().dealloc(NonNull::new_unchecked(ptr), layout)
11    }
12 }
```

## 4.2.2 物理页帧管理器

### ► 接口抽象

在 NoAxiom 中，我们将除了内核 elf 映射区域以外的所有区域均作为物理页帧空间进行管理。物理页帧（Physical Frame）固定分配大小为 4KB，分配时将附带对应物理地址区域的初始化，并支持对于统一物理页帧进行引用计数。物理页帧参与了程序加载、页表分配、I/O 缓冲区等多种功能实现，是 NoAxiom 内核中最重要的内存分配方式之一。

```

1 trait FrameAllocator {
2     fn new() -> Self;
3     fn alloc(&mut self) -> Option<PhysPageNum>;
4     fn dealloc(&mut self, ppn: PhysPageNum);
5 }
```

### ► 具体定义

页帧分配器中维护了一个可变的页帧范围，其中 `current` 代表当前分配到的页帧号，`start` 和 `end` 分别表示页帧区间的起始页号与终止页号。此外还有 `recycled` 数组维护空闲页帧，`frame_map` 用于维护页帧计数映射。对于其具体的 `alloc` 和 `dealloc`，我们使用简易的栈式分配策略，将空闲的物理页帧插入到 `Vec` 当中进行维护。

```

1 pub struct StackFrameAllocator {
2     start: usize,
3     current: usize,
4     end: usize,
5     recycled: Vec<usize>,
6     frame_map: BTreeMap<usize, Weak<FrameTrackerInner>>,
7 }
8 impl FrameAllocator for StackFrameAllocator {
9     fn alloc(&mut self) -> Option<PhysPageNum> {
10         if let Some(ppn) = self.recycled.pop() {
11             Some(ppn.into())
12         } else if self.current == self.end {
13             None
14         } else {
15             self.current += 1;
16             Some((self.current - 1).into())
17         }
18     }
19     fn dealloc(&mut self, ppn: PhysPageNum) {
20         self.recycled.push(ppn.0);
21     }
}
```

### ► 引用计数

为了便于进行写时复制的维护，我们对物理页帧使用 `frame_map: BTeeMap` 进行引用计数维护，其内部保存了物理页帧号到具体的物理页帧追踪器的映射，使用 `Arc` 相关的 `strong_count()` 进行引用计数查询，能够便利地查询当前物理页的引用情况。

```

1 pub fn frame_refcount(ppn: PhysPageNum) -> usize {
2     FRAME_ALLOCATOR
3         .lock()
4         .frame_map
5         .get(&ppn.0)
6         .map_or(0, |x| x.strong_count())
7 }
```

### 4.2.3 资源占用维护

#### ► 物理页帧分配限制

在进行 Unixbench 的 SHELL16 测试点时，我们注意到就算是在 qemu 下运行系统，也容易因为大量申请进程空间而导致物理页帧不足，引发系统崩溃。为了防止这一情况发生，我们尝试在进行大批量的资源分配前进行内存空间的预先检查。例如，我们在进行 elf 文件映射前，会预先查看映射所需空间大小，如果发现映射后系统资源紧张则不予执行。

```

1 pub fn can_alloc(&self, req_num: usize) -> bool {
2     self.stat_remain() >= req_num
3 }
4 pub fn can_alloc_loosely(&self, req_num: usize) -> bool {
5     // reserve 20% more frames for later use
6     self.can_alloc(req_num + req_num / 5 + self.stat_total() / 20)
7 }
8 pub async fn map_elf(...) {
9     if !can_frame_alloc_loosely(frame_req_num) {
10         return_errno!(Errno::ENOMEM, "no enough frames to load elf");
11     }
12 }
```

### 4.3 地址空间维护

在 NoAxiom 中，地址空间被定义为 `MemorySet` 对象保存在任务控制块中，并受到共享自旋锁保护。

`MemorySet` 的原始定义如下所示：

```

1 pub struct MemorySet {
2     pub page_table: SyncUnsafeCell<PageTable>,
3     pub areas: Vec<MapArea>,
4     pub stack: MapArea,
5     pub brk: BrkAreaInfo,
6     pub mmap_manager: MmapManager,
7     pub shm: ShmInfo,
8 }
```

通过对于 `MemorySet` 进行访问，我们可以进行内存映射 `mmap` 维护、程序间断点 `brk` 维护、共享内存 `shm` 维护。在 `NoAxiom` 中，我们使用自旋锁 `Arc<SpinLock<MemorySet>>` 来保证地址空间信息维护的一致性。

## 4.4 延迟分配技术

### 4.4.1 写时复制

在以 `clone(2)` 为代表的进程复制操作过程当中，假如调用时未设置 `VM_CLONE` 标志位，则进程在创建新的子进程时需要将原有的内存信息完整地拷贝一份。通常这一行为将耗费大量的内存空间，为我们的内核带来极大的内存开销负担。

为了尽量避免进程复制时的内存复制开销，我们在内核中引入了写时复制（copy-on-write, COW）技术。具体而言，我们在物理页的页表项标志位（Page Table Entry Flags, PTEFlags）中，额外维护了一个特殊的扩展标记位 COW 用于标记某一页正处于写时复制状态，并取消了对应页面的写标记位（Write Flag）用于捕获对应的用户态写入操作。当一个页面被标记为 COW 时，表示其应当复制，但当前仍由多个进程的地址空间共享。当该页面上发生用户写操作的时候，由于未设置 W 位，该写入操作将被操作系统捕获，并进行实际的页面复制行为。

#### ► 标记位定义

我们在 RV64 和 LA64 的 `PTEFlags` 中均根据实际情况定义了该标记位 `COW`，并在上层的 `MappingFlags` 中创建了对应的虚拟标记位 `cow`，用于上层内核进行访问。

```
1 pub struct MappingFlags: u64 {
2     const COW = bit!(8);
3 }
4 /// RV64_PTE
5 pub struct PTEFlags: u64 {
6     const COW = bit!(8);
7 }
8 /// LA64_PTE
9 pub struct PTEFlags: usize {
10     const COW = bit!(9);
11 }
```

而在内核当中，我们会通过调用 `flags_switch_to_cow` 等函数，实现 COW 页面与普通页面的互转。

此外，值得一提的是，在 LA64 中 W 位的判断方式与其他常规架构存在较大差别。LA64 中，TLB 内部并不存储对应的 W 位，也不会进行相应的合法性检查。因此，假如直接置空 W 位并标记位 COW，硬件仍然会正常地访问该页面而不会触发任何异常。因此在维护 LA64 架构下的 COW 行为时，我们选择了使用置空 V 位的方式来实现页表访问的异常捕获。

```
1 #[inline(always)]
2 pub fn flags_switch_to_cow(flags: &MappingFlags) -> MappingFlags {
3     *flags | MappingFlags::NV | MappingFlags::COW
4 }
5 #[inline(always)]
6 pub fn flags_switch_to_rw(flags: &MappingFlags) -> MappingFlags {
7     *flags & !MappingFlags::NV | MappingFlags::V | MappingFlags::W
8 }
```

## ► 延迟分配

当用户或内核态访问某一携带 COW 标记的页面时，我们将通过各种方式检测到该标记，并通过 `validate` 函数调用 `realloc_cow(vpn, pte)`，最终实现某一页的写时复制分配行为。

在实际分配的时候，我们还会对于页面的引用计数进行检查。当发现某一页的引用计数为 1 时，说明该页面仅由当前线程独享，那么也就没有必要进行写时复制的再次分配了。我们只会将当前页面重置为正常页面，而不会进行额外页面的申请。

```
1 pub fn realloc_cow(&mut self, vpn: VirtPageNum, pte: &PageTableEntry) -> SysResult<()>
2 {
3     let old_ppn = PhysPageNum::from(pte.ppn());
4     let old_flags = pte.flags();
5     let new_flags = flags_switch_to_rw(&old_flags);
6     if frame_refcount(old_ppn) == 1 {
7         self.page_table().set_flags(vpn, new_flags);
8     } else {
9         let frame = frame_alloc().unwrap();
10        let new_ppn = frame.ppn();
11        let mut target = None;
12        // ... iterate memory_set and find the target
13        self.page_table()
14            .remap_cow(vpn, new_ppn, old_ppn, new_flags);
15    }
16    Arch::tlb_flush();
17    Ok(())
}
```

### 4.4.2 懒分配

我们在内存映射（`mmap`）、程序间断点（`brk`）、用户栈（`stack`）的分配过程中使用了懒分配技术。具体而言，在用户程序申请对应资源的时候，我们并不会在实际页表上进行页面信息的映射，而是将对应页面保存到当前 `MemorySet` 中。当用户读取对应页面的数据，会因为页面未映射而触发缺页异常，从而使得访存行为被内存捕获。此时再进行真正的分配过程。

### 4.4.3 合法性检查与分配

我们使用 `validate()` 函数对于某个具体的访存相关异常进行验证，假如检测到位于写时复制或者懒分配区域，则访问地址空间并对于具体的区域进行分配操作。

- 写时复制：我们会首先检查当前发生缺页异常的页面是否真实存在，如果检测到页面存在，则说明该页面可能是一个被标记为 COW 的页面。紧接着检查当前页面的 COW 标志位，假如的确存在，则进行写时复制处理，否则报出地址错误，由对应系统调用或者例外处理例程进行响应。
- 懒分配：假如检测到页面不存在，则说明该页面可能是一个懒分配页面。我们将依次访问内存空间当中维护的 `stack`、`brk`、`mmap` 区域，检查是否的确存在一个对应的页面尚未分配。若未分配，则进行懒分配操作，否则报出地址错误，由对应系统调用或者例外处理例程进行响应。

```

1 pub async fn validate(
2     memory_set: &Arc<SpinLock<MemorySet>>,
3     vpn: VirtPageNum,
4     trap_type: Option<TrapType>,
5     pte: Option<&mut PageTableEntry>,
6 ) -> SysResult<()> {
7     if let Some(pte) = pte {
8         let flags = pte.flags();
9         if flags.contains(MappingFlags::COW) {
10             memory_set.lock().realloc_cow(vpn, pte)?; // copy-on-write allocation
11             Ok(())
12         } else {
13             Err(Errno::EFAULT)
14         }
15     } else {
16         let mut ms = memory_set.lock();
17         if ms.stack.vpn_range.is_in_range(vpn) {
18             let task = current_task().unwrap();
19             ms.lazy_alloc_stack(vpn); // stack lazy allocation
20             Ok(())
21         } else if ms.brk.area.vpn_range.is_in_range(vpn) {
22             ms.lazy_alloc_brk(vpn); // brk(heap) lazy allocation
23             Ok(())
24         } else if ms.mmap_manager.is_in_space(vpn) {
25             drop(ms);
26             lazy_alloc_mmap(memory_set, vpn).await?; // mmap lazy allocation
27             Ok(())
28         } else {
29             Err(Errno::EFAULT)
30         }
31     }
32 }
```

## 4.5 用户指针

在大量系统调用中都涉及用户态传入指针及其合法性检查的步骤。为了规范化用户态指针验证步骤，我们设计了用户指针 `UserPtr` 专门用于进行用户地址空间指针的访问。其本体只包装了一个裸指针地址，其安全性由访问时的合法性检查保证。

```

1 pub struct UserPtr<T = u8> {
2     _phantom: PhantomData<T>,
3     addr: usize,
4 }
```

### 4.5.1 基于异常的指针预检查

在验证用户指针合法性时，我们参考了优秀开源项目 Phoenix<sup>8</sup>的实现思路，设计了适用于 NoAxiom 的指针验证方式。在进行用户指针检查时，arch 层将通过临时替换异常入口并尝试进行一个字节的读写

<sup>8</sup><https://gitlab.eduxiji.net/edugc-group-22026-2376550/T202418123993075-1053>

操作来验证页面的合法性。该设计有效利用了当前 TLB 内缓存的页表映射信息，避免了直接访存导致的 cache miss，有效提升了验证时性能。

由于这一实现与具体架构中异常替换以及异常种类转化相关，我们将具体的指针验证过程放置在了 arch 层，而内核当中则通过 `Arch::check_read(addr)` 和 `Arch::check_write(addr)` 进行指针合法性检查，这两个函数会对于对应地址的读写属性进行检查，若检测出页面异常则会通过 `Err(TrapType)` 返回异常类型。

```
1 pub type UserPtrResult = Result<(), TrapType>;
2 pub trait ArchTrap {
3     fn check_read(addr: usize) -> UserPtrResult;
4     fn check_write(addr: usize) -> UserPtrResult;
5 }
```

就 arch 而言，以 RV64 架构为例，其具体处理 `check_read` 的过程如下，对于读取操作，内核尝试读取一字节的数据，并有可能在这个过程中触发异常。这里需要注意，应当提前关闭 RV64 的压缩指令优化，以防止后续通过 `epc` 自增跳过指令对齐地址情况。

```
1 unsafe fn bare_read(ptr: usize) {
2     asm!(
3         ".option push"
4         ".option norvc"
5         "lb a0, 0(a0) # exception happens here"
6         ".option pop",
7         "in(\"a0\") ptr",
8     );
9 }
```

为了方便实现，我们在全局设置了 per-cpu 的 `USER_PTR_TRAP_TYPE`，在访问用户指针之前将其设置为 `None`，假如访问后检测到其数值发生了变化，则说明访问过程中发生了异常。

```
1 // NoAxiom/lib/arch/src/rv64/trap.rs
2 #[repr(align(64))]
3 #[derive(Copy, Clone, Debug, PartialEq, Eq)]
4 struct Wrapper(TrapType);
5 static mut USER_PTR_TRAP_TYPE: [Wrapper; CPU_NUM] = [Wrapper(TrapType::None); CPU_NUM];
6 unsafe fn before_user_ptr() {
7     RV64::disable_interrupt();
8     set_ptr_entry();
9     // clear trap_type for trap handler
10    USER_PTR_TRAP_TYPE[RV64::get_hartid()] = Wrapper(TrapType::None);
11 }
12 unsafe fn after_user_ptr() -> UserPtrResult {
13     // reload trap_type to check the validity
14     let trap_type = volatile_load(&USER_PTR_TRAP_TYPE[RV64::get_hartid()]).0;
15     let res = match trap_type {
16         TrapType::None => Ok(()),
17         _ => Err(trap_type),
18     };
19     set_kernel_trap_entry();
20     RV64::enable_interrupt();
21     res
22 }
```

当发生异常时，内核会首先跳转到预先设置的异常处理入口 `__kernel_user_ptr_vec` 处理函数中，紧接着将对应的 `USER_PTR_TRAP_TYPE` 设置为从 `scause` 中读取的内容，最后返回主控制流，由主控制流再进行访存，并检测异常类型。

```
1 pub unsafe extern "C" fn kernel_user_ptr_handler() {
2     let hartid = RV64::get_hartid();
3     let scause = scause::read();
4     let stval = stval::read();
5     let sepc = sepc::read();
6     sepc::write(sepc + 4); // skip read
7     USER_PTR_TRAP_TYPE[hartid] = Wrapper(get_trap_type(scause, stval));
8 }
9 unsafe fn check_read(ptr: usize) -> UserPtrResult {
10     before_user_ptr();
11     bare_read(ptr);
12     after_user_ptr()
13 }
14 /// NoAxiom/lib/arch/src/rv64/trap.S
15 __kernel_user_ptr_vec:
16 KERNEL_SAVE_REG
17 call    kernel_user_ptr_handler
18 KERNEL_LOAD_REG
19 sret
```

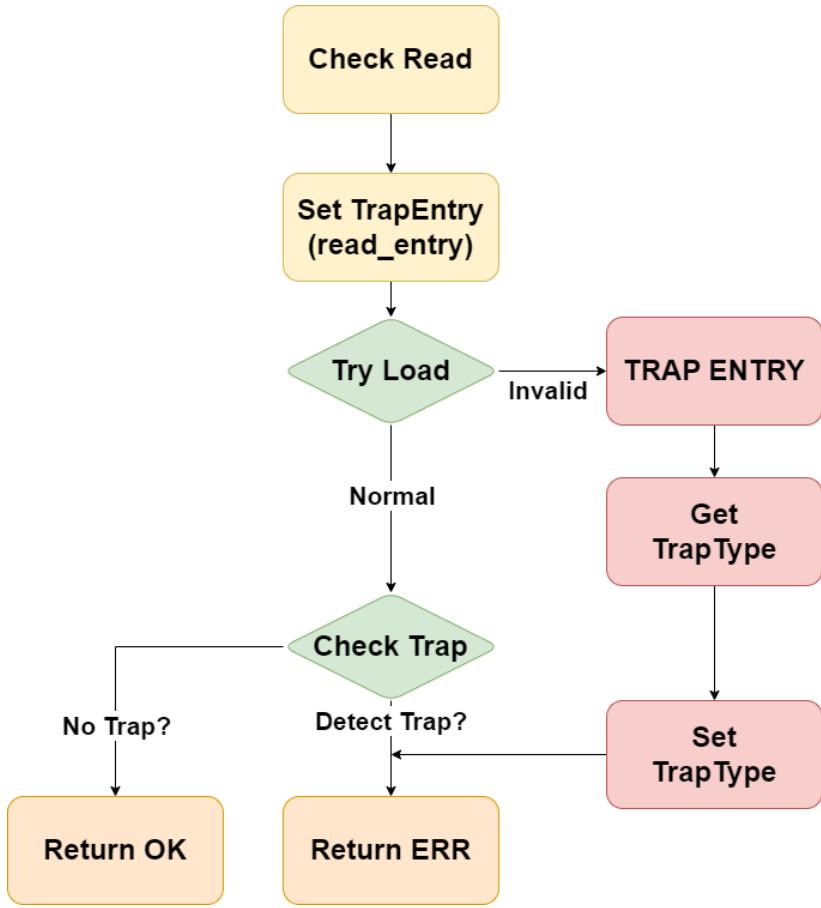


图 12: 基于异常的指针读权限验证

#### 4.5.2 用户指针的延迟分配与异步读取

我们认为检测出用户态指针异常是一件小概率事件，因此在异常发生后的延迟分配过程中，我们无需过度关注性能表现。根据异常的类型不同，我们会访问内存空间并依次对于写时复制、堆栈空间、mmap 空间等进行合法性验证。一旦检测到存在延迟分配行为，则真正进行资源分配。这一逻辑在 `validate()` 函数当中实现。

NoAxiom 系统最大的特点是尽可能追求一切操作的完全异步，在无栈协程架构内尽可能避免让权行为的发生。对于用户指针的懒分配过程，我们也做到了这一点。

由于懒分配中可能存在文件读取行为，该分配过程是有概率进行让权的。我们观察到往年的作品在文件懒分配里的处理都是直接使用阻塞逻辑进行读取，这造成了性能的大量损失。而我们通过将用户指针验证包装为异步函数的方式，实现了用户指针的**异步读取让权**行为。

具体而言，我们以对单个地址的读取为例，详细解释某个指针的合法性验证过程。

当尝试对某一地址进行读取时，首先我们会对于指针的合法性进行检查。这主要包含了两方面：指针是否为空，以及指针对应的物理页是否合法。

对于指针对应页面合法性的检查，我们使用了上一节 4.5.1 中的检查方式，使用 `check_read(addr)` 进行某一单一指针的合法性检查。

此外，当我们检测到当前验证函数执行于某一自旋锁内部的时候，为了防止进程切出导致的死锁，我们选择进行阻塞式的内存验证行为。

```
1 pub async fn read(&self) -> SysResult<T>
2 where
3     T: Copy,
4 {
5     if self.is_null() {
6         return Err(Errno::EFAULT);
7     }
8     match Arch::check_read(self.addr()) {
9         Ok(_) => Ok(unsafe { self.read_unchecked() }),
10        Err(trap_type) => match trap_type {
11            TrapType::LoadPageFault(addr) | TrapType::StorePageFault(addr) => {
12                let task = current_task().unwrap();
13                if check_no_lock() {
14                    task.memory_validate(addr, Some(trap_type), false).await?;
15                } else {
16                    block_on(task.memory_validate(addr, Some(trap_type), true))?;
17                }
18                Ok(unsafe { self.read_unchecked() })
19            }
20            _ => return Err(Errno::EFAULT),
21        },
22    }
23 }
```

例如，当内核试图从用户空间中读取时，将会调用 `ptr.read().await?`，保证用户空间的正确性。其中使用了 Rust 的 `?` 关键字用于简便返回 `Result` 类型返回值中 `Err(_)` 的情况。

```
1 // sys_setitimer
2 let new_value = UserPtr::<ITimerVal>::new(new_value);
3 let new_value = new_value.read().await?;
```

## 五、指令集架构抽象层

### 5.1 双架构支持概览

由于 NoAxiom 的内核架构是从零自主研发的，为了适配赛事的双架构环境，我们设计了独属于 NoAxiom 操作系统内核的硬件抽象层 NoAxiom-Arch 指令集架构抽象层，通过 `Arch` 对象抽象实现了内核与硬件的完全解耦。

目前，NoAxiom 的硬件抽象层支持从 RISCV64、LoongArch64 双架构，支持双架构下从 qemu 启动，支持双架构统一内存管理、TLB 管理，支持双架构统一中断异常处理过程，此外还支持时钟维护等其他架构相关的解耦维护。

NoAxiom-Arch 向上仅暴露了 `Arch` 对象及其相关实现的若干 `trait` 属性。我们应用面向对象多态的设计思想，通过 `trait` 重载虚函数的方式，将硬件抽象接口统一为一致的接口，供给内核调用。

### 5.2 硬件架构解耦设计

NoAxiom 中，所有涉及指令集架构的操作均被统一为 `Arch` 对象中的不同方法供 kernel 调用。编译时，我们将根据选择的架构不同，为 `Arch` 选定对应的底层实现。

```
1 #[cfg(target_arch = "loongarch64")]
2 pub type Arch = la64::LA64;
3 #[cfg(target_arch = "riscv64")]
4 pub type Arch = rv64::RV64;
```

在 NoAxiom-Arch 中，硬件相关的不同功能被具体分类为中断使能控制、内嵌汇编、中断异常陷入处理、时间管理、系统启动、架构信息等子模块，这些子模块被包含在对应的特性 `trait` 当中，并由顶层特性 `FullArch` 统一管理。我们规定 `Arch` 应当实现 `FullArch`，并通过该特性向内核传递具体的方法调用。

```
1 pub trait ArchFull:
2     ArchInt + ArchAsm + ArchSbi + ArchTrap + ArchTime + ArchBoot + ArchInfo {}
3 pub struct LA64;
4 impl crate::ArchFull for LA64 {}
5 pub struct RV64;
6 impl crate::ArchFull for RV64 {}
```

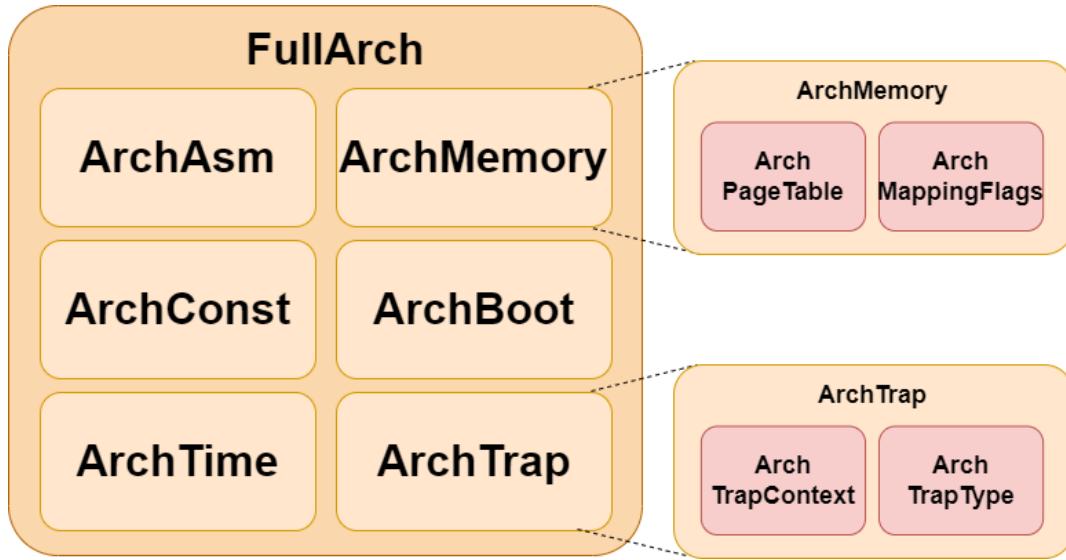


图 13: 架构层整体结构

### 5.3 硬件初始化

在 NoAxiom 中，硬件相关的初始化分为两部分：

1. 系统从 bootloader 启动后刚进入内核代码时，HAL 层通过 `entry` 汇编代码块，进行初始的内核相关控制寄存器配置。我们称之为 `early_init`。
2. 系统从 HAL 层进入到 kernel 层进行初始化后，在初始化过程中需要额外进行与硬件相关的初始化过程。我们称之为 `kernel_arch_init`。

#### 5.3.1 启动时初始化

由于汇编代码要求程序给定 `_entry`，我们并没有为 `entry` 相关的函数配置对应的 `trait` 内函数，而是直接嵌入到代码中。

当程序从 bootloader 跳转到 `_entry` 进行初始化的时候，我们会进行启动时临时页表配置，这主要是为了适配我们内核的高低半地址空间的设计。除此之外，还会进行内核栈空间的配置，我们通过预留的内核栈空间指针，将不同核心绑定到对应的内核栈上。

对于两个架构而言，我们都会预先为每个核心预留一个内核栈，并在初始化的时候进行内核栈的设置。下文中的汇编代码片段展示了 rv64 中的内核栈初始化过程，该过程在 la 中基本一致。具体而言，在启动时，我们会获取当前核心的 cpuid 号，并通过该 id 号与单一内核栈大小相乘，以此获取对应的内核栈顶地址。

```

1 #[link_section = ".bss.kstack"]
2 pub(crate) static BOOT_STACK: [u8; KERNEL_STACK_SIZE * CPU_NUM] = [0; KERNEL_STACK_SIZE
   * CPU_NUM];
3
4 asm! (
5     add      t0, a0, 1

```

```

6     slli    t0, t0, {kernel_stack_shift}
7     la      sp, {boot_stack}
8     add    sp, sp, t0
9     li      s0, {kernel_addr_offset}
10    or      sp, sp, s0
11    li      t1, {kernel_addr_offset}
12    or      gp, gp, t1
13  ");

```

## ►RISC-V64 临时页表初始化

rv64 中，由于分配了内核位于高半地址空间但内核空间并未分配，我们需要使用预先定义的初始化时页表来充当临时页表进行翻译。该页表会将 0xffff\_fc00\_8000\_0000 开头的地址映射到低半空间的物理地址上，实现内核空间的正常访问。

```

1 #[link_section = ".data.prepage"]
2 static PAGE_TABLE: [usize; PTE_PER_PAGE] = {
3     let mut arr: [usize; PTE_PER_PAGE] = [0; PTE_PER_PAGE];
4     macro_rules! page_table_config {
5         ($($id:expr, $addr:expr)*) => {
6             $($arr[$id] = ($addr << 10) | 0xcf;)*
7         };
8     }
9     page_table_config! {
10        2, 0x80000
11        0x102, 0x80000
12    };
13    arr
14 };
15
16 // activate page table
17 asm!(
18     la      t0, {page_table}
19     srli    t0, t0, 12
20     li      t1, 8 << 60
21     or      t0, t0, t1
22     csrw    satp, t0
23     sfence.vma
24 );

```

## ►LoongArch64 内存映射窗口配置

在 LA64 中，我们为了便于维护，选择直接配置 LA64 下的内存映射窗口，将 0x8000 与 0x9000 开头的地址段标记为内核可访问，以此实现高地址空间的内存映射。

```

1 asm!(
2     ori      $t0, $zero, 0x1          # CSR_DMW1_PLV0
3     lu52i.d $t0, $t0, -2048        # UC, PLV0, 0x8000 xxxx xxxx xxxx
4     csrwr   $t0, 0x180            # LOONGARCH_CSR_DMWIN0
5     ori      $t0, $zero, 0x11        # CSR_DMW1_MAT | CSR_DMW1_PLV0
6     lu52i.d $t0, $t0, -1792        # CA, PLV0, 0x9000 xxxx xxxx xxxx
7     csrwr   $t0, 0x181            # LOONGARCH_CSR_DMWIN1
8 )

```

## ► 内核入口跳转

在全部初始化完成之后，我们将会跳转到内核当中预留的真正入口 `_boot_hart_init` 中，退出汇编进行进一步的初始化过程。下面展示了 rv64 中相关的跳转操作。

```
1  asm! ("
2      mv      a1, gp
3      la      t0, {entry}
4      or      t0, t0, t1
5      mv      a0, tp
6      jalr    t0
7  "),
8  entry = sym _boot_hart_init
```

### 5.3.2 内核功能初始化

在 NoAxiom 中，我们设计了 `trait ArchBoot` 用于进行系统初始化相关的内核调用接口，其内部搭载了 `arch_init` 函数用于进行初始化。

```
1  pub trait ArchBoot {
2      fn arch_init();
3 }
```

#### ►LoongArch64

对于 la64 而言，其初始化过程包含了浮点寄存器、异常中断向量入口、时钟模块、页表翻译 (tlb) 相关的初始化。其具体初始化过程不再赘述。

```
1  impl ArchBoot for LA64 {
2      fn arch_init() {
3          freg_init();
4          trap_init();
5          time_init();
6          tlb_init();
7      }
8 }
```

#### ►RISC-V64

而对于 rv64 而言，由于页表翻译过程由硬件接管，我们在一开始启动的时候并不需要对于 tlb 相关寄存器进行特殊配置。此外，由于 RISCV 特有的 SBI 接口的存在，我们不需要特地配置时钟相关的初始化。只需要进行例外入口和浮点寄存器初始化即可。

```
1  impl ArchBoot for RV64 {
2      fn arch_init() {
3          trap_init();
4          freg_init();
5      }
6 }
```

## 5.4 页表管理

### 5.4.1 页表抽象

在我们的 `ArchMemory` 中，对于页表抽象进行了给定。

```

1 /// memory management arch trait
2 pub trait ArchMemory {
3     const PHYS_MEMORY_START: usize;
4     const PHYS_MEMORY_END: usize;
5     const KERNEL_ADDR_OFFSET: usize;
6     type PageTable: ArchPageTable;
7     fn tlb_init();
8     fn tlb_flush();
9     fn current_root_ppn() -> usize;
10    fn activate(ppn: usize, is_kernel: bool);
11 }

```

可以看到，我们设计了 `ArchPageTable` 特性，并定义了 `type PageTable: ArchPageTable` 用于向内核提供统一的页表定义接口，借此限定底层的具体实现。

而对于 `ArchPageTable` 的具体实现，我们进一步提供了 `PageTableEntry` 接口，用于限定单个页表项的访问行为。并提供了获取当前根页表页帧号、创建新页表、激活某一页表等一系列功能，向内核提供统一的访问接口。

```

1 pub trait ArchPageTable {
2     type PageTableEntry: ArchPageTableEntry;
3
4     /// virtual address width
5     const VA_WIDTH: usize;
6     /// index level number
7     const INDEX_LEVELS: usize;
8
9     fn root_ppn(&self) -> usize;
10    fn new(root_ppn: usize) -> Self;
11    fn activate(&self);
12 }

```

#### 5.4.2 页表项抽象

在 RV64 和 LA64 中，页表项的维护策略并不完全相同，根据我们对于硬件架构的了解，我们将页表项异同总结如下：

- 相同点：均保持了物理页帧号 + 标记位的结构，且均具备 RWV 等基本的标记位。
- 不同点：物理页帧号保存的偏移不同，标记位的具体位置不同。

我们将相同点统合为 `ArchPageTableEntry` 这一类型，该类型通过 `MappingFlags` 向内核提供一个软件模拟的页表控制位，并在与硬件层交互时，将其转化为具体的对应 PTE 中的控制位。同时，由于 PTE 中的物理页帧号的维护需求基本相同，我们选择使用统一传入的物理页帧号，与 `MappingFlags` 拼接，最终形成对应硬件架构下的 PTE 值。

```

1 pub trait ArchPageTableEntry: Into<usize> + From<usize> + Clone {
2     /// create a new page table entry from ppn and flags
3     fn new(ppn: usize, flags: MappingFlags) -> Self;
4     /// get the physical page number

```

```

5     fn ppn(&self) -> usize;
6     /// get the pte permission flags
7     fn flags(&self) -> MappingFlags;
8     /// set flags
9     fn set_flags(&mut self, flags: MappingFlags);
10    /// clear all data
11    fn reset(&mut self);
12    /// is valid dir
13    fn is_allocated(&self) -> bool;
14
}

```

这里需要强调，`ArchPageTableEntry` 向内核暴露的接口中，传入的控制参数均为 `MappingFlags` 类型，但我们在底层设置时，将该 flag 进一步转化为硬件定义的 `PTEFlags` 再进行设置。

我们对于虚拟页表项控制位 `MappingFlags` 的定义如下所示。可以看到，除了常规的 `V, R, W, X` 等标记位以外，我们定义了三个特殊的标记位 `COW, NV, PT`。其中，`COW` 位用于写时复制标记，`NV` 用于强制标识当前页面不合法，`PT` 用于强制标识当前页面为页表中的非叶节点，主要用于简化 LA 下的页表遍历过程。

```

1 bitflags! {
2     /// Mapping flags for page table.
3     #[derive(Debug, Clone, Copy, PartialEq, Eq, PartialOrd, Ord, Hash)]
4     pub struct MappingFlags: u64 {
5         const V = bit!(0);
6         const R = bit!(1);
7         const W = bit!(2);
8         const X = bit!(3);
9         const U = bit!(4);
10        const G = bit!(5);
11        const A = bit!(6);
12        const D = bit!(7);
13        const COW = bit!(8);
14        /// OS-specific: virt bit to specify the invalid page
15        const NV = bit!(62);
16        /// OS-specific: virt bit for page table
17        const PT = bit!(63);
18    }
19 }

```

## ►RV64 页表项转化

在 RV64 架构下，PTE 中保存的控制位的真实定义如下，可以看到，与 `MappingFlags` 并不完全相同，我们还需要进一步的映射来实现真实控制位的设置。

```

1 bitflags! {
2     #[derive(Debug, Clone, Copy, PartialEq, Eq, PartialOrd, Ord, Hash)]
3     pub struct PTEFlags: u64 {
4         const V = bit!(0);
5         const R = bit!(1);
6         const W = bit!(2);
7         const X = bit!(3);
8         const U = bit!(4);
9         const G = bit!(5);
10        const A = bit!(6);
11        const D = bit!(7);
12
}

```

```

12     const COW = bit!(8);
13 }
14 }
```

而在进行 `MappingFlags` 到 `PTEFlags` 的转化中，我们将实现对应控制位的一一映射行为。在下面的代码中，我们通过宏定义实现了 `RWXUGAD` 等控制位的转换。

```

1 impl From<MappingFlags> for PTEFlags {
2     fn from(flags: MappingFlags) -> Self {
3         if flags.is_empty() {
4             Self::empty()
5         } else if flags.contains(MappingFlags::PT) {
6             Self::V
7         } else {
8             let mut res = Self::empty();
9             macro_rules! set {
10                 ($($flag:ident),*) => {
11                     $(
12                         if flags.contains(MappingFlags::$flag) {
13                             res |= PTEFlags::$flag;
14                         }
15                     )*
16                 };
17             }
18             if flags.contains(MappingFlags::V) && !flags.contains(MappingFlags::NV) {
19                 res |= PTEFlags::V;
20             }
21             set!(R, W, X, U, G, A, D, COW);
22             res
23         }
24     }
25 }
```

## ►LA64 页表项转化

LA64 下的页表项与 RV64 类似，不过由于出现了更多未在 `MappingFlags` 中描述的控制位，我们对于未定义位使用了统一的映射方式进行控制。

这里值得一提的，为了优化 LA 下的 TLB 访问行为，我们选择将所有页表内非叶节点的项的 flag 置空，这样就可以避免 TLB 查询时对于控制位的额外计算了；但在 RV 中，中间项也必须设置 V 标记位，用于指示当前页表项的合法性。下面列出了 RV 和 LA 下对于 `MappingFlags::PT` 的不同处理。

```

1 // RV64: set V flag
2 if flags.contains(MappingFlags::PT) {
3     Self::V
4 }
5 // LA64: set empty, for TLB iteration pref
6 if flags.contains(MappingFlags::PT) {
7     return PTEFlags::empty();
8 }
```

### 5.4.3 TLB 维护

LA64 与 RV64 下的 TLB 重填维护存在很大不同，其主要原因是 LA64 下的 TLB 重填过程由软件负责，而 RV64 下该行为则几乎完全依赖硬件实现。这导致 LA64 下需要额外地进行 TLB 重填相关操作的维护。

#### ►TLB 初始化

具体而言，LA64 下在系统初始化时，我们就会进行 TLB 重填的相关寄存器初始化。我们的内核的页面大小固定为 4KB，因此会在初始化时向对应 csr 寄存器固定设置页面大小为 PS\_4K。除此之外，还会进行 pwcl 和 pwch 寄存器的配置，用于控制页表翻译过程中每级在虚地址中对应的数位范围。

最终当所有寄存器初始化完成后，我们会调用 `set_tlb_refill_entry` 设置 TLB 重填异常的入口，并冲刷 TLB 使能该设置。

```
1 pub fn tlb_init_inner() {
2     // Page Size 4KB
3     const PS_4K: usize = 0x0c;
4     tlbidx::set_ps(PS_4K);
5     stlbps::set_ps(PS_4K);
6     tlbrehi::set_ps(PS_4K);
7     // Set Page table entry width
8     pwcl::set_pte_width(8);
9     // Set Page table width and offset
10    pwcl::set_ptbase(12);
11    pwcl::set_ptwidth(9);
12    pwcl::set_dir1_base(21);
13    pwcl::set_dir1_width(9);
14    pwcl::set_dir2_base(30);
15    pwcl::set_dir2_width(9);
16    pwch::set_dir3_base(39);
17    pwch::set_dir3_width(9);
18    set_tlb_refill_entry(__tlb_refill as usize);
19 }
```

#### ►TLB 重填异常

我们在 `__tlb_refill` 当中使用汇编嵌入了一段 TLB 重填的处理例程，在陷入 TLB 重填异常时，将会通过 `tlbreentry` 寄存器的值跳转到该函数对应的 PC 值，并进行 TLB 重填处理。重填过程将会遍历四级页表并将获取到的 PTE 填入 TLB 中。

```
1 .globl __tlb_refill
2 .balign 4096
3 __tlb_refill:
4     csrwr    $t0, LA_CSR_TLBRSAVE
5     csrrd    $t0, LA_CSR_PGD
6     lddir    $t0, $t0, 3
7     lddir    $t0, $t0, 2
8     lddir    $t0, $t0, 1
9     ldpte    $t0, 0
10    ldpte   $t0, 1
11    tlbfill
12    csrrd    $t0, LA_CSR_TLBRSAVE
13    ertn
```

## 5.5 异常处理

### 5.5.1 异常处理抽象

异常处理过程由于涉及到特定 CSR 寄存器的读写操作，因此也与硬件架构强耦合。我们选择使用 `ArchTrap` 特性实现异常处理的抽象接口。该接口包含了异常初始化 `trap_init`、上下文切换 `trap_restore`、异常状态读取 `read_trap_type` 等操作，同时对于用户指针合法性验证，使用 `check_read` 和 `check_write` 函数进行维护。

```
1 pub trait ArchTrap {
2     type TrapContext: ArchTrapContext;
3     fn trap_init();
4     fn trap_restore(_cx: &mut <Self as ArchTrap>::TrapContext);
5     fn read_epc() -> usize;
6     fn read_trap_type(cx: Option<&mut <Self as ArchTrap>::TrapContext>) -> TrapType;
7     fn check_read(addr: usize) -> UserPtrResult;
8     fn check_write(addr: usize) -> UserPtrResult;
9 }
```

### 5.5.2 异常类型抽象

由于不同架构下对于异常类型的定义不相同，我们设计了统一的抽象异常类型 `TrapType` 用于向内核提供统一的异常访问接口。定义如下所示，其中与访存相关的错误会额外包含一个 `usize` 表示的地址，用于指示导致异常的地址，用于内核验证内存时使用。

```
1 pub enum TrapType {
2     Breakpoint,
3     SysCall,
4     Timer,
5     Unknown,
6     SupervisorExternal,
7     SupervisorSoft,
8     StorePageFault(usize),
9     LoadPageFault(usize),
10    InstructionPageFault(usize),
11    IllegalInstruction(usize),
12    None,
13 }
```

### 5.5.3 异常类型转化

#### ►RV64 下的异常类型转化

我们的异常类型主要基于 RV64 进行设计，因此 RV64 架构下的异常类型转化接近于直接映射的行为。具体而言，我们会预先读取 RV64 下用于指示异常的 `scause` 与 `stval` 的值，并将其传入 `get_trap_type` 中，由该函数进行进一步的映射。

```
1 pub fn get_trap_type(scause: Scause, stval: usize) -> TrapType {
2     match scause.cause() {
3         Trap::Exception(Exception::LoadFault) => TrapType::Unknown,
```

```

4     Trap::Exception(Exception::UserEnvCall) => TrapType::SysCall,
5     Trap::Interrupt(Interrupt::SupervisorTimer) => TrapType::Timer,
6     Trap::Exception(Exception::StorePageFault) => TrapType::StorePageFault(stval),
7     Trap::Exception(Exception::StoreFault) => TrapType::StorePageFault(stval),
8     Trap::Exception(Exception::InstructionPageFault) => TrapType::
9         InstructionPageFault(stval),
10    Trap::Exception(Exception::IllegalInstruction) => TrapType::IllegalInstruction(
11        stval),
12    Trap::Exception(Exception::LoadPageFault) => TrapType::LoadPageFault(stval),
13    Trap::Interrupt(Interrupt::SupervisorExternal) => TrapType::SupervisorExternal,
14    Trap::Interrupt(Interrupt::SupervisorSoft) => TrapType::SupervisorSoft,
15    _ => panic!("unknown trap type: {:?}", scause.cause()),
}
}

```

## ►LA64 下的异常类型转化

LA64 下的异常类型转化相对较为复杂，由于 LA64 下定义了 `AddressNotAligned` 等需要特殊处理的异常类型，我们无法建立其与 `TrapType` 的直接映射，而是需要先进行合法性检查与处理再进行返回。例如，我们对于 `AddressNotAligned` 的处理依赖于 LA 本身的架构相关指令，该异常的处理发生在 `get_trap_type` 函数内部，而返回给内核的将是 `TrapType::None`，用于指示已经预先完成了异常处理。

```

1 fn get_trap_type(tf: Option<&mut TrapContext>) -> TrapType {
2     let estat = estat::read();
3     let badv = badv::read().vaddr();
4     match estat.cause() {
5         Trap::Exception(e) => {
6             match e {
7                 Exception::Breakpoint => TrapType::Breakpoint,
8                 Exception::AddressNotAligned => {
9                     unsafe { emulate_load_store_insn(tf.unwrap()) }
10                    TrapType::None
11                }
12                Exception::Syscall => TrapType::SysCall,
13                Exception::StorePageFault | Exception::PageModifyFault => TrapType::
14                    StorePageFault(badv),
15                    // ... 下略
16                }
17            }
18        }
}

```

## 5.6 上下文维护

### 5.6.1 寄存器抽象

我们的内核在系统调用、设置栈帧等操作的时候，需要对于已经保存的上下文中的特定寄存器进行操作。但不同架构下，寄存器的编号并不一致，因此我们设计了 `TrapArgs` 用于提供统一的上下文中寄存器访问接口。

具体而言，我们提供了如下的寄存器接口供内核访问。

表 6: 任务控制块成员定义

名称	含义	对应内核功能
EPC	异常 PC 值	异常处理
RA	返回地址	信号处理等
SP	栈帧地址	信号处理、线程创建等
RES	系统调用返回值	系统调用
A0~A5	系统调用参数	系统调用
SYSCALL	系统调用号	系统调用
TLS	线程指针	线程创建

### 5.6.2 上下文抽象

我们在 ArchTrap 特性中给定了类型 TrapContext，该类型用于向内核提供上下文接口。而其内部实现了 ArchTrapContext 特性，用于衔接底层具体实现，其定义如下：

```

1 pub trait ArchTrapContext:
2     Index<TrapArgs, Output = usize> + IndexMut<TrapArgs, Output = usize>
3 {
4     type FloatContext: ArchUserFloatContext;
5     fn app_init(cx: &mut Self, entry: usize, sp: usize) -> &Self;
6     fn freg_mut(&mut self) -> &mut Self::FloatContext;
7     fn gprs(&self) -> &[usize; 32];
8     fn gprs_mut(&mut self) -> &mut [usize; 32];
9     fn get_syscall_id(&self) -> usize;
10    fn get_syscall_args(&self) -> [usize; 6];
11 }
```

可以看到，我们为 ArchTrapContext 实现了按照 TrapArgs 进行索引的特性，借此内核就可以通过如 cx[TrapArgs::SP] 的方式访问到某一上下文的具体寄存器了。我们以 RV64 为例，其实现如下所示。而在 LA64 中也有对应的寄存器映射方式，具体不再赘述。

```

1 impl Index<TrapArgs> for TrapContext {
2     type Output = usize;
3     fn index(&self, index: TrapArgs) -> &Self::Output {
4         match index {
5             TrapArgs::EPC => &self.sepc,
6             TrapArgs::RA => &self.x[reg_id::RA],
7             TrapArgs::SP => &self.x[reg_id::SP],
8             TrapArgs::RES => &self.x[reg_id::A0],
9             TrapArgs::A0 => &self.x[reg_id::A0],
10            TrapArgs::A1 => &self.x[reg_id::A1],
11            TrapArgs::A2 => &self.x[reg_id::A2],
12            TrapArgs::A3 => &self.x[reg_id::A3],
13            TrapArgs::A4 => &self.x[reg_id::A4],
14            TrapArgs::A5 => &self.x[reg_id::A5],
15            TrapArgs::TLS => &self.x[reg_id::TP],
16            TrapArgs::SYSCALL => &self.x[reg_id::A7],
17        }
18    }
19 }
```

而在内核中想要修改上下文寄存器值时，我们将通过 `cx[TrapArgs]` 具体寄存器值的访问。以系统调用为例，我们在处理系统调用时大量访问了上下文中的寄存器，包括 `EPC`、`RES`，并通过调用 `get_syscall_id()` 和 `get_syscall_args()` 实现了系统调用参数的获取。

```
1 impl Task {
2     pub async fn syscall(self: &Arc<Self>, cx: &mut TrapContext) -> isize {
3         cx[TrapArgs::EPC] += 4;
4         let res = get_syscall_result(
5             Syscall::new(self)
6                 .syscall(cx.get_syscall_id(), cx.get_syscall_args())
7                 .await,
8         );
9         cx[TrapArgs::RES] = res as usize;
10        clear_current_syscall();
11        res
12    }
13 }
```

此外，信号相关的上下文 `uc_mcontext` 也与架构强相关，因此我们添加了 `SigContext` 类型用于架构相关的上下文维护。

## 六、设备抽象层

### 6.1 设备抽象设计

#### 6.1.1 设备抽象层概览

设备抽象层是 NoAxiom 系统 HAL 层中的重要中间层，负责统合底层的架构抽象层与基础内存管理模块，并向内核提供统一的设备管理接口与设备驱动接口。

设备抽象层的设计主要由设备探测、设备管理与设备驱动三个子模块构成。设备探测模块负责在机器启动时从设备树和总线上获取设备信息，确定合法设备的基址。设备管理模块负责注册管理设备，并向内核提供统一的设备获取接口。设备驱动模块根据探测到的设备信息，提供具体设备的底层驱动实现。

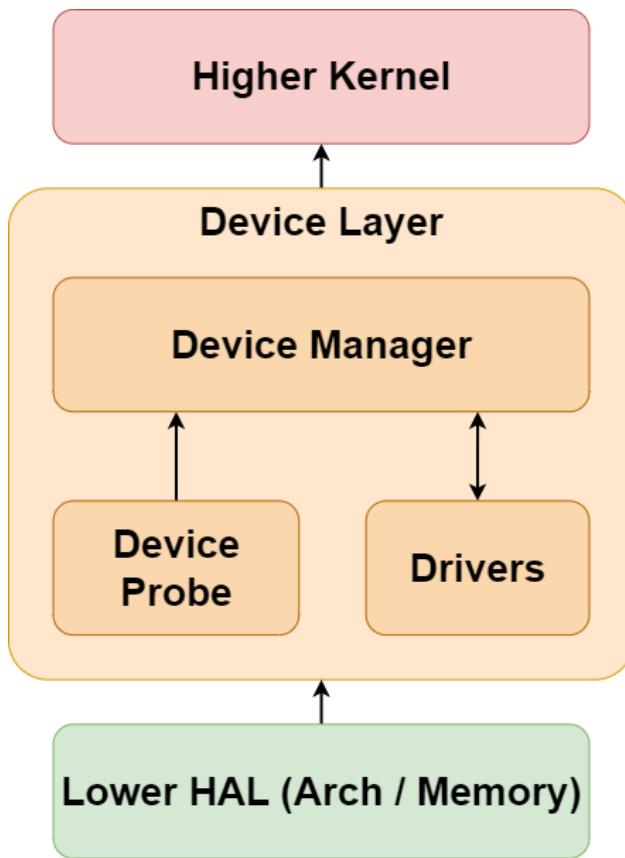


图 14: 设备抽象层架构

#### 6.1.2 设备类型抽象

为了便于设备管理，我们定义了统一的基础设备类型 `DeviceType` 用于表示设备类型。

```
1 pub enum DeviceType {  
2     Block(BlockDeviceType),  
3     Net(NetDeviceType),
```

```

4     Display(DisplayDeviceType),
5     Interrupt(InterruptDeviceType),
6     Char(CharDeviceType),
7     Power(PowerDeviceType),
8     Kernel,
9     Probe,
10    Unknown,
11 }

```

对于具体设备，例如块设备，我们使用 `BlockDeviceType` 来表示其使用的具体驱动类型。对于块设备而言，可能使用三种类型的驱动：virtio，visionfive 板卡上的 sdcard 驱动，以及 ls2k1000 板卡上的 ahci 驱动。我们通过该枚举类型指示其驱动注册，对于待定的驱动则使用 `Unknown` 进行表示。

```

1 pub enum BlockDeviceType {
2     Virtio,
3     VF2Sdcard,
4     LS2k1000Ahci,
5     Unknown,
6 }

```

### 6.1.3 设备驱动抽象

为了统一管理设备驱动，我们基于 Rust 中的面向对象思想，使用 trait 特性进行设备驱动的维护。

我们使用统一的根 trait `Device` 用于表示基础的设备驱动类型，并对于不同的具体设备，定义不同的驱动特性。例如，块设备继承自 `Device`，在 `Device` 基础上，还需要实现 `read()`; `write()` 两种方法以提供读写接口。

```

1 pub trait Device: Sync {}
2 pub trait BlockDevice: Send + Sync + Device + InterruptDevice { /* ... */ }

```

## 6.2 设备探测

### 6.2.1 设备树探测

NoAxiom 内核在初始化基础架构与内存管理模块后，就会执行设备驱动的初始化，其中设备树（Device Tree Blob）的基地址会被传递到驱动层进行解析。我们会解析出当前运行环境下的 `MMIO` 总线基地址、`PCI` 总线基地址，以及 `PLIC`，其中因为 Loongarch64 架构本身不支持平台级中断控制器，所以在 `DtbInfo` 中采用了预编译进行控制。

对于 NoAxiom 内核而言，设备树信息是由 bootloader 或内核直接嵌入给定，该信息在内核运行过程中静态不变。因此我们定义了全局的设备树基址 `DTB_ADDR`，解析时根据运行的平台差异，选择是否调用 `DTB_ADDR` 完成设备树信息的解析与保存。

```

1 static DTB_ADDR: Once<usize> = Once::new();
2 pub fn get_dtb_initializer() -> DtbInitializerType {
3     #[cfg(feature = "qemu")]
4     {
5         match crate::probe::dtb::get_dtb_addr() {
6             Some(addr) => DtbInitializerType::Ptr(addr),

```

```

7         None => panic!("DTB address not initialized"),
8     }
9 }
10 #[cfg(not(feature = "qemu"))]
11 {
12     const FDT: &[u8] = include_bytes!("./jh7110-visionfive-v2.dtb");
13     DtbInitializerType::Ref(FDT)
14 }
15 }
```

获取设备树基址后，设备层将从设备树基址中解析各个设备信息，并通过 `DeviceConfig` 将设备基址、设备类型等信息注册到全局设备信息注册器 `DeviceConfigManager` 中。

```

1 pub struct DeviceConfig {
2     pub dev_type: DeviceType,
3     pub region: MmioRegion,
4     pub conf_type: DeviceConfigType,
5 }
6 pub struct DeviceConfigManager {
7     pub devices: Vec<DeviceConfig>,
8 }
```

为了便于进行设备注册，我们创建了全局设备注册表 `OF_INITIALIZERS`。而其用于注册的信息则来自各个设备的 `DeviceTreeInfo` 特性，该特性会向设备提供 `OF_TYPE` 等设备树相关信息，用于设备匹配与转换。通过该注册表，NoAxiom 能够快捷地提取设备中预先通过特性 `DeviceTreeInfo` 给定的设备信息并进行设备注册。

```

1 pub trait DeviceTreeInfo: Device {
2     const OF_TYPE: &'static str;
3     const DEVICE_TYPE: &'static DeviceType;
4     const DEVICE_CONFIG_TYPE: &'static DeviceConfigType = &DeviceConfigType::DeviceTree
5         ;
6 }
7 macro_rules! device {
8     ($device:ty) => {
9         (
10             <$device>::OF_TYPE,
11             <$device>::DEVICE_TYPE,
12             <$device>::DEVICE_CONFIG_TYPE,
13         )
14     };
15 }
16 pub(super) const OF_INITIALIZERS: &[(&str, &DeviceType, &DeviceConfigType)] = &[
17     device!(PlicDevice),
18     device!(VF2SdcardDevice),
19     device!(LsAhciDevice),
20 ];
```

`DeviceConfig` 中保存的设备信息此后再经由实例化函数 `device_realize` 进行实际的设备初始化。在该函数中，我们遍历已经从设备树中获取到的设备信息，并根据不同的设备类型进行不同的初始化操作。

```

1 fn dtb_realize() {
2     let manager = DEV_CONFIG_MANAGER.get().unwrap();
```

```

3     for config in manager.devices.iter() {
4         match config.conf_type {
5             DeviceConfigType::VirtioMmio => virtio_mmio_realize(config),
6             DeviceConfigType::PciEcam => pci_realize(config),
7             DeviceConfigType::DeviceTree => of_realize(config),
8         }
9     }
10 }
```

### 6.2.2 总线设备嗅探

因为在比赛背景下，由于不同架构下 QEMU 参数的不同（RISC-V64 架构下是 MMIO 总线上挂载的块设备，而 Loongarch64 架构下是在 PCI 总线上挂载的块设备），所以我们区别于以往队伍设备探测时采用固定 MMIO 基址的方式，借鉴了 ArceOS<sup>9</sup>的嗅探方法，实现了**多总线下的设备自动探测**，成功的将设备驱动的注册与底层运行平台实现了解耦。具体实现上，因为设备驱动也是对于内核运行而言静态不变的，所以使用了全局结构 `ALL_DEVICES` 用来保存所有内核所需的设备驱动。除此之外，我们还为网络模块实现了 `LoopBack` 回环设备驱动。

```

1 pub fn probe_bus() {
2     ALL_DEVICES.as_ref_mut().probe_mmiobus_devices().ok();
3     ALL_DEVICES.as_ref_mut().probe_pcibus_devices().ok();
4     ALL_DEVICES.as_ref_mut().add_net_device(LoopBackDev::new());
5     debug!("[driver] probe succeed");
6 }
```

### 6.2.3 PCI 总线下的设备嗅探

首先，利用 `DTB` 里保存的 PCI 总线基址，获得**根节点**(Root Complex)，这是 PCI 总线的根节点，为我们后续提供与 PCI 总线的管理和通信。

```

1 pub(crate) fn probe_pcibus_devices(&mut self) -> DevResult<()> {
2     let base_vaddr = dtb_info().pci_ecam_base | arch::Arch::KERNEL_ADDR_OFFSET;
3     let mut root = unsafe { PciRoot::new(base_vaddr as *mut u8, Cam::Ecam) };
4     ...
5 }
```

在这之后，我们会遍历 `PCI_BUS_END` 个总线，并且遍历每个总线上所有的设备信息，得到挂载在上面的设备 `BDF(Bus, Device, Function)` 信息。

```

1 /// An identifier for a PCI bus, device and function.
2 #[derive(Copy, Clone, Debug, Eq, PartialEq)]
3 pub struct DeviceFunction {
4     /// The PCI bus number, between 0 and 255.
5     pub bus: u8,
6     /// The device number on the bus, between 0 and 31.
7     pub device: u8,
8     /// The function number of the device, between 0 and 7.
9     pub function: u8,
```

<sup>9</sup><https://github.com/arceos-org/arceos/tree/main/modules/axdriver/src/bus>

10 }

然后，我们便对 PCI 设备尝试进行配置，如果 BAR(Base Address Register) 寄存器配置无误的话，就可以注册相应的驱动了。

```

1 pub(crate) fn probe_pcibus_devices(&mut self) -> DevResult<()> {
2     ...
3     for bus in 0..=PCI_BUS_END as u8 {
4         for (bdf, dev_info) in root.enumerate_bus(bus) {
5             if dev_info.header_type != HeaderType::Standard {
6                 continue;
7             }
8             match config_pci_device(&mut root, bdf, &mut allocator) {
9                 Ok(_) => {
10                     // Register driver
11                     ...
12                 }
13                 Err(e) => log::warn!(
14                     "failed to enable PCI device at {}({}): {:?}", 
15                     bdf,
16                     dev_info,
17                     e
18                 ),
19             }
20         }
21     }
22     Ok(())
23 }
```

```

1 fn config_pci_device(
2     root: &mut PciRoot,
3     bdf: DeviceFunction,
4     allocator: &mut Option<PciRangeAllocator>,
5 ) -> DevResult<()> {
6     // Check bar info
7     ...
8
9     // Enable the device.
10    let (_status, cmd) = root.get_status_command(bdf);
11    root.set_command(
12        bdf,
13        cmd | Command::IO_SPACE | Command::MEMORY_SPACE | Command::BUS_MASTER,
14    );
15    Ok(())
16 }
```

#### 6.2.4 MMIO 总线下的设备嗅探

相较于 PCI 总线，MMIO(Memory-Mapped I/O) 总线下的嗅探就相对简明很多，它将外围设备直接映射到内存空间，可以通过普通访存指令直接交互。在设备树读取了 MMIO 信息后，利用 virtio-driver-async(我们自主 fork 修改的异步驱动库，后面会有详细介绍) 库直接进行注册。目前我们注册了块设备(Block Device)，未来会加入更多的字符设备注册。

```

1 pub fn probe_mmiobus_devices(&mut self) -> DevResult<()> {
2     ...
3     for (addr, _size) in &dtb_info().virtio_mmio_regions {
4         // Register driver
5     }
6     Ok(())
7 }
```

## 6.3 驱动层

### 6.3.1 驱动层结构

我们参考 VFS 虚拟文件系统的思想，使用了 `trait` 抽象出多种不同具体物理设备的一致功能，向内核提供了清晰简明的接口。

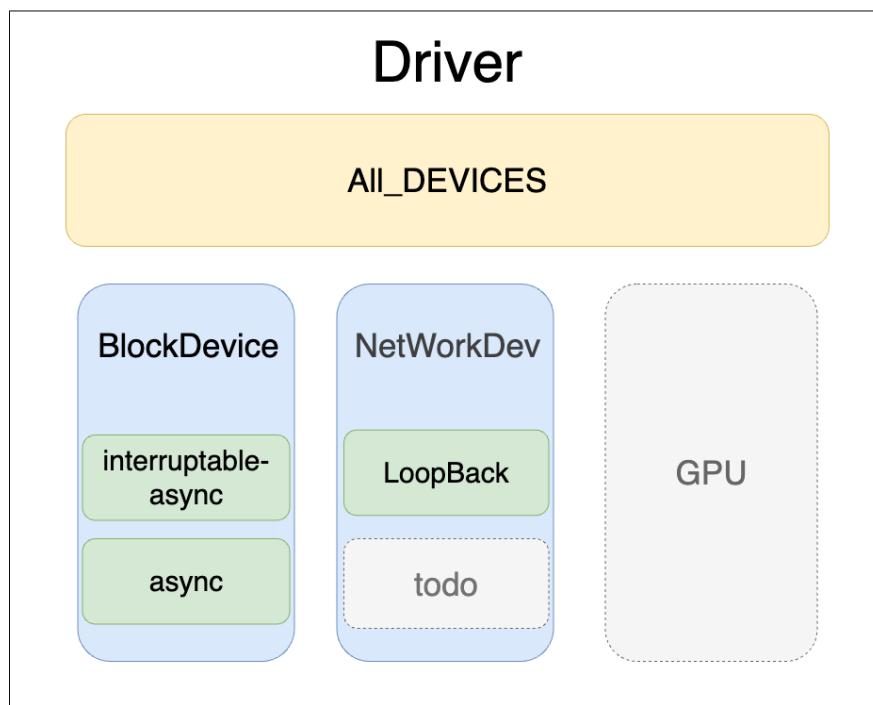


图 15: 驱动层结构

对于块设备而言，需要实现 `BlockDevice trait`：

```

1 pub trait BlockDevice: Send + Sync {
2     fn device_name(&self) -> &'static str;
3     fn handle_interrupt(&self) -> DevResult<()> {
4         unimplemented!("{} not implement handle_interrupt!", self.device_name())
5     }
6     async fn read(&self, id: usize, buf: &mut [u8]) -> DevResult<usize> {
7         // the buf's len is multiple of BLOCK_SIZE
8         unimplemented!("{} not implement read!", self.device_name())
9     }
}
```

```

10     async fn write(&self, id: usize, buf: &[u8]) -> DevResult<usize> {
11         // the buf's len is multiple of BLOCK_SIZE
12         unimplemented!("{} not implement read!", self.device_name())
13     }
14     async fn sync_all(&self) -> DevResult<()> {
15         unimplemented!("{} not implement sync_all!", self.device_name())
16     }
17 }

```

特别的，我们为块设备驱动 `trait` 中添加了 `handle_interrupt` 方法，用来处理外部中断。在我们的异步块设备驱动中注册并使用。

对于网卡设备，需要实现 `NetWorkDev trait`：

```

1 pub trait NetWorkDev: Send + Sync {
2     /// get the MAC address of the network card
3     fn mac(&self) -> EthernetAddress;
4     fn iface_name(&self) -> String;
5     /// get the network card ID
6     fn nic_id(&self) -> usize;
7     fn poll(&self, sockets: &mut SocketSet) -> DevResult<()>;
8     /// get the smoltcp interface type
9     fn inner_iface(&self) -> &SpinLock<Interface>;
10 }

```

未来如果需要添加更多的设备驱动（如 GPU），只需要完成对应设备的 `trait` 约束即可。这些驱动在经过了设备嗅探之后，会保存在全局结构 `ALL_DEVICES` 里，供上层使用。

```

1 // ALL_DEVICES struct
2
3 pub struct Devices {
4     net: Option<NetDevice>,
5     blk: Option<BlkDevice>,
6     display: Option<DisplayDevice>,
7     // ... more devices
8 }

```

## 6.4 异步块设备驱动

为充分利用 NoAxiom 的无栈协程架构，我们创新性地设计了异步块设备驱动<sup>10</sup>。我们将底层块设备忙等的时间利用内核提供的 `yield_now` 方法进行让权，进而减少了 CPU 忙等的时间，提高了 I/O 操作的吞吐率。

具体而言，我们修改了 `add_notify_wait_pop` 函数，库原本函数的逻辑是忙等可用环：

```

1 pub fn add_notify_wait_pop<'a>(...) -> Result<u32> {
2     ...
3     // Wait until there is at least one element in the used ring.
4     while !self.can_pop() {
5         spin_loop();
6     }
7     ...
8 }

```

<sup>10</sup><https://github.com/YuXuaann/virtio-drivers-async>

经过修改后，我们把 `spin_loop` 改为了 `yield`。

```
1 pub async fn add_notify_wait_pop<'a>(...) -> Result<u32> {
2     ...
3     // Wait until there is at least one element in the used ring.
4     loop {
5         if self.can_pop() {
6             break;
7         }
8         yield_now().await;
9     }
10    ...
11 }
```

最终使得 `virtio-driver-async` 库对外暴露的 `read_blocks` 函数成为了异步函数，而且采用的是轮询的异步逻辑，最终完美适配于驱动层定义的 `BlockDevice trait`。而且此种方案因为是基于对 `virtio-driver` 库的 fork 更改，它可以适用于 MMIO 总线和 PCI 总线下的块设备，所以不管是对于 RISC-V64 还是 Loongarch64 架构，我们都可以使用异步驱动。

特别的，对于 RISC-V64 架构而言，因其有对外部中断良好支持的中断控制器 (PLIC)，于是我们也采用了往届优秀队伍设计的基于中断的异步块设备驱动<sup>11</sup>，并将其运用到我们的内核里。

对于原本的 `virtio-driver` 驱动，它完成一次 I/O 操作的过程是这样的：

1. 从描述符表中申请三个空闲描述符，填写 `req`、`resp`、`buf` 字段，其中 `buf` 字段是写操作需要的
2. 通过写特定 MMIO 寄存器的方式通知设备进行 I/O 读写
3. 设备完成读写操作后，会把已完成的 I/O 的描述符放到已用环中
4. 驱动通过轮询去查看已用环中是否有需要的描述符，如果有就返回结果

以下是原 `virtio` 驱动的读扇区方法：

```
1 impl<const N: usize> VirtIOBlock<N> {
2     pub fn read_sector(&self, block_id: usize, buf: &mut [u8]) -> Result<()> {
3         ...
4         q.add_buf(&[req.as_buf()], &[buf, resp.as_buf_mut()])
5             .expect("[virtio] virtual queue add buf error");
6         h.notify(0);
7         while !q.can_pop() {} // busy loop here
8         q.pop_used()?;
9         match resp.status {
10             BlockRespStatus::Ok => Ok(()),
11             _ => Err(VirtIOError::IOError),
12         }
13     }
14 }
```

但是因为这样轮询的效率比较差，所以我们利用设备的中断功能，结合无栈协程机制，实现了唤醒机制。具体实现上是在 `VirtIOBlock` 结构加上一个 `Event` 成员，用来监听 I/O 操作完成的事件。

<sup>11</sup><https://github.com/HUST-OS/tornado-os>

```

1 impl<const N: usize> VirtIOBlock<N> {
2     pub async fn read_sector_event(&self, sector_id: usize, buf: &mut [u8]) -> Result
3         <()> {
4             ...
5             q.add_buf(&[req.as_buf()], &[buf, resp.as_buf_mut()])
6                 .expect("[virtio] virtual queue add buf error");
7             h.notify(0);
8             listener.await; // wait for wake
9             q.pop_used()?;
10            match resp.status {
11                BlockRespStatus::Ok => Ok(()),
12                _ => Err(VirtIOError::IOError),
13            }
14        }
}

```

其中 `wake` 当且仅当接收到外部中断时发生，利用前述 `BlockDevice trait`，完美实现中断处理函数：

```

1 #[async_trait]
2 impl BlockDevice for VirtIOAsyncBlock {
3     fn handle_interrupt(&self) -> DevResult<()> {
4         unsafe {
5             self.handle_interrupt()
6                 .expect("virtio handle interrupt error!");
7             self.0.wake_ops.notify(1); // wake the listener
8         };
9         Ok(())
10    }
11 }

```

## 6.5 PLIC 中断控制器

### 6.5.1 初始化

为适配于需要中断的设备驱动，我们会在 RISC-V64 架构利用 PLIC 库<sup>12</sup>配置 `PLIC`。在设备树初始化完成后，会将 `PLIC` 基址存在 `DTB_INFO` 结构体中，供 `PLIC` 初始化使用。类似于 `DTB_INFO`，`PLIC` 的结构同样采用 `Once` 包裹，在初始化时调用 `PLIC.call_once`，供未来中断处理使用。

```

1 pub static PLIC: Once<PLIC<CPU_NUM>> = Once::new();
2
3 pub fn init() {
4     let plic_addr = dtb_info().plic | arch::Arch::KERNEL_ADDR_OFFSET;
5     debug!("PLIC addr: {:#x}", plic_addr);
6     let privileges = [2; CPU_NUM];
7     let plic = PLIC::new(plic_addr, privileges);
8     PLIC.call_once(|| plic);
9
10    ...
11 }

```

初始化过程中还会设置各种设备中断的优先级，目前仅设置了块设备中断（中断号为 1）的优先级：

<sup>12</sup><https://github.com/os-module/plic>

```

1 pub fn init() {
2     ...
3     let priority = match () {
4         #[cfg(feature = "interruptable_async")]
5         () => 1,
6         #[cfg(feature = "async")]
7         () => 0,
8     };
9     let irq = 1;
10    let plic = PLIC.get().unwrap();
11    plic.set_priority(irq, priority);
12    ...
13 }

```

除此之外，PLIC 会为所有 CPU 逻辑核心配置每个核心在内核态 (Supervisor) 的接受中断的阈值 (该核仅接受优先级大于等于该阈值的中断)。因为**中断委托**机制的存在，所以在内核态发生的中断我们也可以在内核态进行处理，这样内核就基本接管了一切外部中断。

```

1 pub fn register_to_hart(hart: u32) {
2     let plic = PLIC.get().unwrap();
3     let irq = 1;
4     plic.enable(hart, Mode::Supervisor, irq);
5     plic.set_threshold(hart, Mode::Supervisor, 0);
6     log::info!("Register irq {} to hart {}", irq, hart);
7 }
8
9 pub fn init() {
10     ...
11     for i in 0..CPU_NUM {
12         register_to_hart(i as u32);
13     }
14 }

```

### 6.5.2 外部中断处理

因为我们内核中断均在内核态处理，所以对外部库提供的 `claim` 和 `complete` 方法做了简易封装：

```

1 pub fn claim() -> u32 {
2     let plic = PLIC.get().unwrap();
3     let hart = Arch::get_hartid();
4     plic.claim(hart as u32, Mode::Supervisor)
5 }
6 pub fn complete(irq: u32) {
7     let plic = PLIC.get().unwrap();
8     let hart = Arch::get_hartid();
9     plic.complete(hart as u32, Mode::Supervisor, irq);
10 }

```

在内核接收到中断后，就会调用 `handle_irq()` 方法，然后就可以根据中断号选择合适的设备进行中断处理，只需要执行 `trait` 中定义的 `handle_interrupt` 方法。目前仅有块设备注册了中断处理。

```

1 pub fn handle_irq() {
2     #[cfg(feature = "interruptable_async")]

```

```

3     {
4         let irq = plic::claim();
5         if irq == 1 {
6             get_blk_dev()
7                 .handle_interrupt()
8                 .expect("handle interrupt error");
9         } else {
10            log::error!("[driver] unhandled irq: {}", irq);
11        }
12        plic::complete(irq);
13    }
14 #[cfg(feature = "async")]
15 {
16     unreachable!("shouldn't accept interrupt!");
17 }
18 }
```

## 6.6 LoopBack 回环设备

利用 smoltcp 库提供的 Loopback 结构，实现了内核的回环设备驱动。

```

1 pub struct LoopBackDev {
2     pub interface: SpinLock<Interface>,
3     pub dev: SpinLock<Loopback>,
4 }
5
6 impl LoopBackDev {
7     pub fn new() -> Self {
8         let mut device = Loopback::new(Medium::Ethernet);
9         let iface = {
10             ...
11         };
12         Self {
13             interface: SpinLock::new(iface),
14             dev: SpinLock::new(device),
15         }
16     }
17 }
```

并为其实现 NetworkDev trait，其中比较关键的是 poll 和 inner\_iface 方法。前者是更新网卡状态的重要方法，后者在 smoltcp 中 TcpSocket 中的 connect 函数中需要使用，

```

1 impl NetworkDev for LoopBackDev {
2     fn mac(&self) -> EthernetAddress {
3         EthernetAddress([0x00, 0x00, 0x00, 0x00, 0x00, 0x00])
4     }
5
6     fn iface_name(&self) -> String {
7         String::from("lo")
8     }
9
10    fn nic_id(&self) -> usize {
11        // loopback's netcard id is 0
12        0
13    }
14 }
```

```
14
15     fn poll(&self, sockets: &mut iface::SocketSet) -> DevResult<()> {
16         let mut iface = self.interface.lock();
17         let mut device_guard = self.dev.lock();
18         let device = device_guard.deref_mut();
19         let res = iface.poll(Instant::from_millis(get_time_ms() as i64), device,
20                             sockets);
21         if res {
22             Ok(())
23         } else {
24             Err(Errno::EAGAIN)
25         }
26     }
27
28     fn inner_iface(&self) -> &SpinLock<Interface> {
29         &self.interface
30     }
}
```

## 七、文件系统

### 7.1 VFS

VFS 是操作系统内核中用于抽象和统一各种文件系统的一个软件层。它提供了一个通用的接口，使得上层应用程序能够以统一的方式访问不同类型的物理文件系统和设备，而无需关心底层的具体实现细节。

#### 7.1.1 总体结构设计

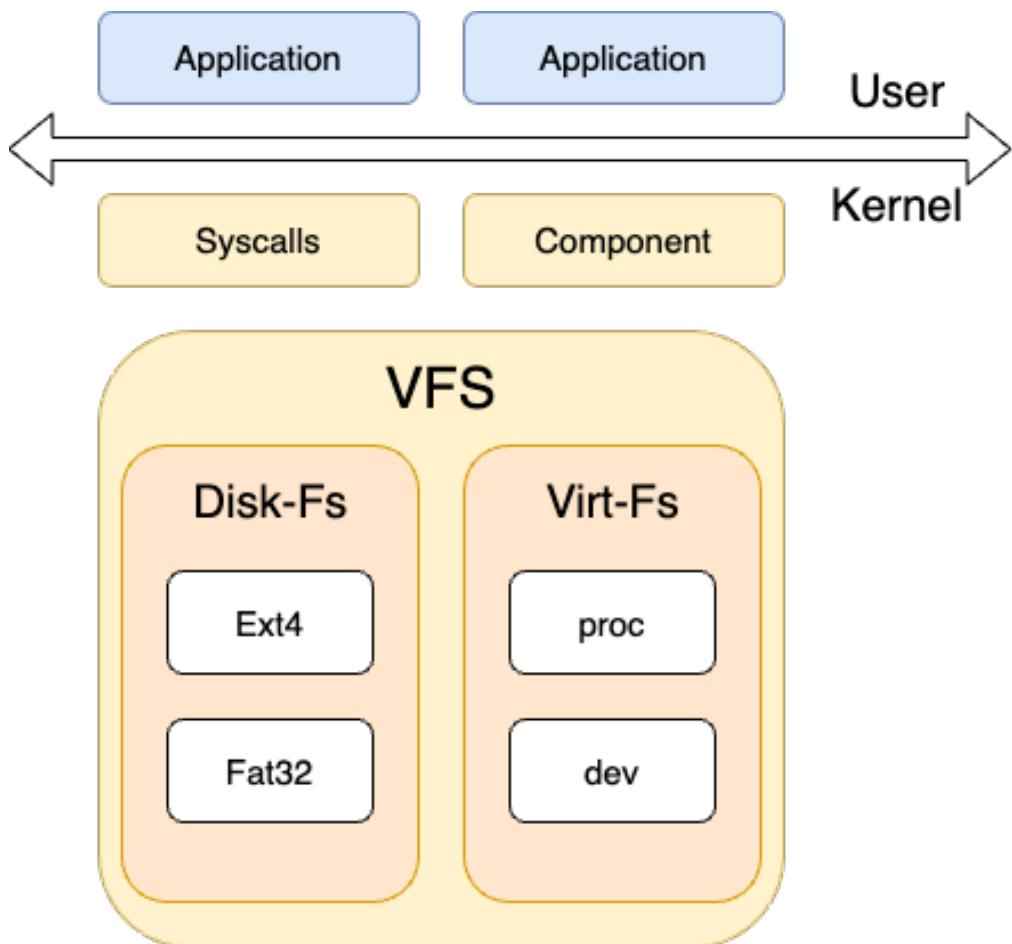


图 16: VFS

#### 7.1.2 目录项 Dentry

目录项是构成 VFS 树型结构的结点结构。它维护了文件系统中所有文件的树形结构联系，是文件名和 [inode](#) 之间的桥梁。它在内存中构成了文件系统的目录树结构，它提供了基于路径的文件的增删查的方法。

目前 Dentry 实现了功能完备、实现清晰的路径解析功能，并且全面支持符号链接的路径解析，同样也支持硬链接功能。

```
1 // NoAxiom/kernel/src/fs/vfs/basic/dentry.rs
2 pub struct DentryMeta {
3     /// The name of the dentry
4     name: String,
5     /// The super block of the dentry
6     pub super_block: Arc<dyn SuperBlock>,
7     /// The parent of the dentry, None if it is root
8     parent: Option<Weak<dyn Dentry>>,
9     /// The children of the dentry
10    children: SpinLock<BTrieMap<String, Arc<dyn Dentry>>>,
11    /// The inode of the dentry, None if it is negative
12    inode: SpinLock<Option<Arc<dyn Inode>>>,
13 }
14
15 #[async_trait]
16 pub trait Dentry: Send + Sync + DowncastSync {
17     /// Get the meta of the dentry
18     fn meta(&self) -> &DentryMeta;
19     /// Open the file associated with the dentry
20     fn open(self: Arc<Self>) -> SysResult<Arc<dyn File>>;
21     /// Get new dentry from name
22     fn from_name(self: Arc<Self>, name: &str) -> Arc<dyn Dentry>;
23     /// Create a new dentry with `name` and `mode`
24     async fn create(self: Arc<Self>, name: &str, mode: InodeMode) -> SysResult<Arc<dyn Dentry>>;
25     /// Get the name of the dentry
26     fn name(&self) -> String {
27         self.meta().name.clone()
28     }
29 }
```

目录项在路径解析时产生，并且永远存在于 VFS 树上，直到文件被删除（或者被重命名）时，才会从 VFS 树上消亡。若一个 Dentry 的 Inode 域为 `None`，称之为负目录项。

```
1 // NoAxiom/kernel/src/fs/vfs/basic/dentry.rs
2 pub fn walk_path(
3     self: &Arc<Self>,
4     task: &Arc<Task>,
5     path: &Vec<&str>,
6 ) -> SysResult<Arc<dyn Dentry>> {
7     Ok(self.__walk_path(Some(task), path, 0, 0)?)
8 }
9
10 fn __walk_path(
11     self: &Arc<Self>,
12     task: Option<&Arc<Task>>,
13     path: &Vec<&str>,
14     step: usize,
15     jumps: usize,
16 ) -> SysResult<(Arc<dyn Dentry>, usize)> {
17     const SYMLINK_MAX_STEP: usize = 12;
18     if unlikely(jumps >= SYMLINK_MAX_STEP) {
19         error!("[walk_path] symlink too many times, jumps: {}", jumps);
```

```

20         return Err(Errno::ELOOP);
21     }
22     if step == path.len() {
23         return Ok((self.clone(), jumps));
24     }
25     if unlikely(self.is_negative()) {
26         error!("[walk_path] {} is negative", self.name());
27         return Err(Errno::ENOENT);
28     }
29
30     let inode = self.inode().expect("should have inode!");
31
32     let entry = path[step];
33     match entry {
34         "." => self.__walk_path(task, path, step + 1, jumps),
35         ".." => {
36             if let Some(parent) = self.parent() {
37                 parent.__walk_path(task, path, step + 1, jumps)
38             } else {
39                 error!("[walk_path] {} has no parent", self.name());
40                 Err(Errno::ENOENT)
41             }
42         }
43         name => {
44             if let Some(symlink_path) = inode.symlink() {
45                 let (tar, new_jumps) = self.__symlink_jump(task, &symlink_path, jumps +
46                     1)?;
47                 return tar.__walk_path(task, path, step, new_jumps);
48             }
49
50             if inode.file_type() != InodeMode::DIR {
51                 error!("[walk_path] {} is not a dir", self.name());
52                 return Err(Errno::ENOTDIR);
53             }
54
55             // Only check search permissions after confirming it's a directory
56             if let Some(task) = task {
57                 if task.user_id().fsuid() != 0 {
58                     if unlikely(!self.can_search(task)) {
59                         error!("[walk_path] has no search access");
60                         return Err(Errno::EACCES);
61                     }
62                 }
63                 if let Some(child) = self.get_child(name) {
64                     return child.__walk_path(task, path, step + 1, jumps);
65                 }
66             match self.clone().open(&FileFlags::empty()) {
67                 Ok(file) => block_on(file.load_dir()).expect("can not load dir!");
68                 Err(_) => return Err(Errno::ENOENT);
69             }
70             if let Some(child) = self.get_child(name) {
71                 return child.__walk_path(task, path, step + 1, jumps);
72             }
73
74             Err(Errno::ENOENT)

```

```
75     }
76 }
77 }
```

区别于大部分没有 Dentry 的 OS 设计，通过 Dentry 可以把路径解析功能解耦出来，并且可以更优雅高效地维护路径解析功能。具体流程如下：

1. 进行一系列合法性 check，包括权限是否支持加载目录、是否为 Dir 类型的目录项、是否为负目录项等等，若有则返回相应错误码；如若到达路径解析的终点，则返回当前 Dentry 结构体；
2. 在当前目录项查找路径中的下一目录项，若存在且并非符号链接文件，则进行递归查找，返回到第一步；
3. 在上一步中，若找不到目录项，则进行 load 操作得到下一目录项，如果仍然找不到则返回 ENOENT；若找到的为符号链接文件，则进行符号链接跳转，然后再递归查找，返回第一步。

下面是符号链接的跳转实现，即为读取 Inode 中存储的符号链接路径，然后再次调用 `__walk_path` 方法进行跳转。其中若符号链接跳转次数超过设定阈值，则会返回 ELOOP。

```
1 pub fn symlink_jump(
2     self: &Arc<Self>,
3     task: &Arc<Task>,
4     symlink_path: &str,
5 ) -> SysResult<Arc<dyn Dentry>> {
6     Ok(self.__symlink_jump(Some(task), symlink_path, 0)? .0)
7 }
8
9 fn __symlink_jump(
10    self: &Arc<Self>,
11    task: Option<&Arc<Task>>,
12    symlink_path: &str,
13    jumps: usize,
14 ) -> SysResult<(Arc<dyn Dentry>, usize)> {
15     let abs = symlink_path.starts_with('/');
16     let components = path::resolve_path(symlink_path)?;
17     let res = if abs {
18         root_dentry().__walk_path(task, &components, 0, jumps)
19     } else {
20         self.parent()
21             .expect("must have parent")
22             .__walk_path(task, &components, 0, jumps)
23     };
24     match res {
25         Ok((dentry, new_jumps)) => {
26             if let Ok(inode) = dentry.inode() {
27                 if let Some(symlink_path) = inode.symlink() {
28                     debug!("[__symlink_jump] Following symlink: {}", symlink_path);
29                     return dentry.__symlink_jump(task, &symlink_path, jumps + 1);
30                 }
31             }
32             Ok((dentry, new_jumps))
33         }
34         Err(e) => Err(e),
35     }
36 }
```

```
35     }
36 }
```

NoAxiom 同样支持在 VFS 层面上的硬链接功能，如若需要在文件系统上持久化保留，只需要具体文件系统支持即可，调用具体文件系统的 `link` 方法即可实现。在用户执行硬链接系统调用时，我们只需要申请一个空 `Dentry`，并且复制一个指向旧文件 `Inode` 的 `Arc` 指针即可，这样利用 rust 原子引用计数指针可以自动实现析构功能（当所有指向 `Inode` 的指针消失时自动析构）。当读取到该硬链接文件时，自然会使用其链接对象的 `Inode` 读取对应的文件内容。

```
1 pub async fn create_link(
2     self: &Arc<Self>,
3     target: Arc<dyn Dentry>,
4     name: &str,
5 ) -> SysResult<isize> {
6     debug_assert!(!target.is_negative());
7     if self.inode()?.file_type() != InodeMode::DIR {
8         error!("[Vfs::link] {} is not a dir", self.name());
9         return Err(errno::ENOTDIR);
10    }
11
12    if let Some(old) = self.get_child(name) {
13        error!(
14            "[Vfs::link] {} already exists in {}",
15            old.name(),
16            self.name()
17        );
18        return Err(errno::EEXIST);
19    }
20
21    let inode = target.inode()?;
22    let son = self.clone().from_name(name);
23    son.set_inode(inode.clone());
24    debug!("[Vfs::linkto] set_inode {} to {}", name, target.name());
25
26    let nlink = inode.meta().inner.lock().nlink;
27    inode.meta().inner.lock().nlink += 1;
28
29    self.add_child(son);
30
31    Ok(nlink as isize)
32 }
```

### 7.1.3 元信息 Inode

元信息记录了文件的一些基本信息，例如文件大小 `size` 以及系统调用所需的 `stat` 和 Loongarch 架构系统调用所需的 `statx` 结构体，并且保存了 `symlink` 结构。它主要提供关于查询文件元信息的一系列方法。

目前已经全面支持符号链接以及 `uid`、`gid` 的维护。

```
1 pub struct InodeMeta {
2     /// The inode id, unique in the file system
```

```

3     pub id: usize,
4     /// The inner data of the inode, maybe modified by multiple tasks
5     pub inner: Mutex<InodeMetaInner>,
6     /// The mode of file
7     pub inode_mode: AtomicU32,
8     /// The super block of the inode
9     pub super_block: Arc<dyn SuperBlock>,
10    /// The page cache of the file, managed by the `Inode`
11    pub page_cache: Option<()>,
12
13    pub uid: AtomicU32,
14    pub gid: AtomicU32,
15    symlink: Mutex<Option<String>>, // for symlink, the target path
16 }
17
18 pub struct InodeMetaInner {
19     /// The number of links to the inode
20     pub nlink: usize,
21     /// The size of the file
22     pub size: usize,
23     /// The state of the file
24     state: InodeState,
25     /// Last access time.
26     pub atime_sec: usize,
27     pub atime_nsec: usize,
28     /// Last modification time.
29     pub mtime_sec: usize,
30     pub mtime_nsec: usize,
31     /// Last status change time.
32     pub ctime_sec: usize,
33     pub ctime_nsec: usize,
34 }
35
36 #[async_trait]
37 pub trait Inode: Send + Sync + DowncastSync {
38     fn meta(&self) -> &InodeMeta;
39     fn stat(&self) -> SysResult<Stat>;
40     async fn truncate(&self, _new: usize) -> SysResult<()> {
41         panic!("this inode not implemented truncate");
42     }
43 }
44
45 impl dyn Inode {
46     pub fn statx(&self, mask: u32) -> SysResult<Statx>
47 }
```

除此之外，`Inode` 还为那些支持缓存功能的文件提供了页缓存 `PageCache`。`Inode` 和文件是一一对应的关系，但当文件被用户 `close` 掉后，考虑到局部性原理，用户很可能再次打开该文件，因此我们仍会保留其 `Inode` 以及相应的缓存内容，进而避免因反复开关文件导致的 `PageCache` 写回开销，这部分内容可见后文缓存设计。

#### 7.1.4 文件 File

文件是进程打开文件在内存中的实例，它提供了一系列关于 I/O 操作（包括 `read/write`, `load`, `poll`）的方法，也实现了通过其 `Inode` 去使用页缓存加速读写的方法。

```
1 pub struct FileMeta {
2     /// File flags, may be modified by multiple tasks
3     flags: SpinLock<FileFlags>,
4     /// The position of the file, may be modified by multiple tasks
5     pub pos: AtomicUsize,
6     /// Pointer to the Dentry
7     dentry: Arc<dyn Dentry>,
8     /// Pointer to the Inode
9     pub inode: Arc<dyn Inode>,
10 }
11
12 #[async_trait]
13 pub trait File: Send + Sync + DowncastSync {
14     /// Get the size of file
15     fn size(&self) -> usize {
16         self.meta().inode.size()
17     }
18     async fn page_cache(&self) -> Option<AsyncMutexGuard<'_, PageCache>> {
19         self.meta().inode.page_cache().await
20     }
21     /// Get the dentry of the file
22     fn dentry(&self) -> Arc<dyn Dentry> {
23         self.meta().dentry.clone()
24     }
25     /// Get the meta of the file
26     fn meta(&self) -> &FileMeta;
27     /// Read data from file at `offset` to `buf`, not for kernel other modules
28     async fn base_read(&self, offset: usize, buf: &mut [u8]) -> SyscallResult;
29     /// Readlink data from file at `offset` to `buf`
30     async fn base_readlink(&self, buf: &mut [u8]) -> SyscallResult;
31     /// Write data to file at `offset` from `buf`, not for kernel other modules
32     async fn base_write(&self, offset: usize, buf: &[u8]) -> SyscallResult;
33     /// Load directory into memory, must be called before read/write explicitly,
34     /// only for directories
35     async fn load_dir(&self) -> Result<(), Errno>;
36     /// Delete dentry, only for directories
37     async fn delete_child(&self, name: &str) -> Result<(), Errno>;
38     /// IOCTL command
39     fn ioctl(&self, cmd: usize, arg: usize) -> SyscallResult;
40     fn poll(&self, req: &PollEvent, waker: Waker) -> PollEvent;
41 }
42
43 impl_downcast!(sync File);
44
45 impl dyn File {
46     pub async fn read_at(&self, offset: usize, buf: &mut [u8]) -> SyscallResult
47     pub async fn write_at(&self, offset: usize, buf: &[u8]) -> SyscallResult
48 }
```

每一个 `File` 结构体都会持有一个指向具体文件系统文件的强引用计数指针，这么做可以保证多个

进程打开同一个文件时不需要消耗过多的内存资源，并且借由 rust 的强引用计数指针管理文件的生命周期，优雅地自动实现文件的析构功能，在所有持有文件指针的进程退出（或这些进程的用户程序调用 `sys_close` 主动丢弃文件指针）后，自动 `Drop`。

我们还利用了原子变量 `AtomicUsize pos`，去标记文件读写指针的位置，避免了使用锁的开销。通过 `store-release` 和 `load-acquire` 的内存顺序进行读写操作，既保证了正确性，又维持了高性能。

### 7.1.5 超级块 SuperBlock

超级块主要提供了文件系统相关的信息。前面所述的结构中，也会存在一个指针指向超级块，以指示当前文件所属的文件系统。

```
1 pub struct SuperBlockMeta {
2     /// The device of the file system, None if it is a virtual file system
3     device: Option<Arc<&'static dyn BlockDevice>>,
4     /// The file system
5     file_system: Arc<dyn FileSystem>,
6     /// The root of the file system, use weak to avoid reference cycle
7     root: Once<Weak<dyn Dentry>>,
8 }
```

## 7.2 缓存设计

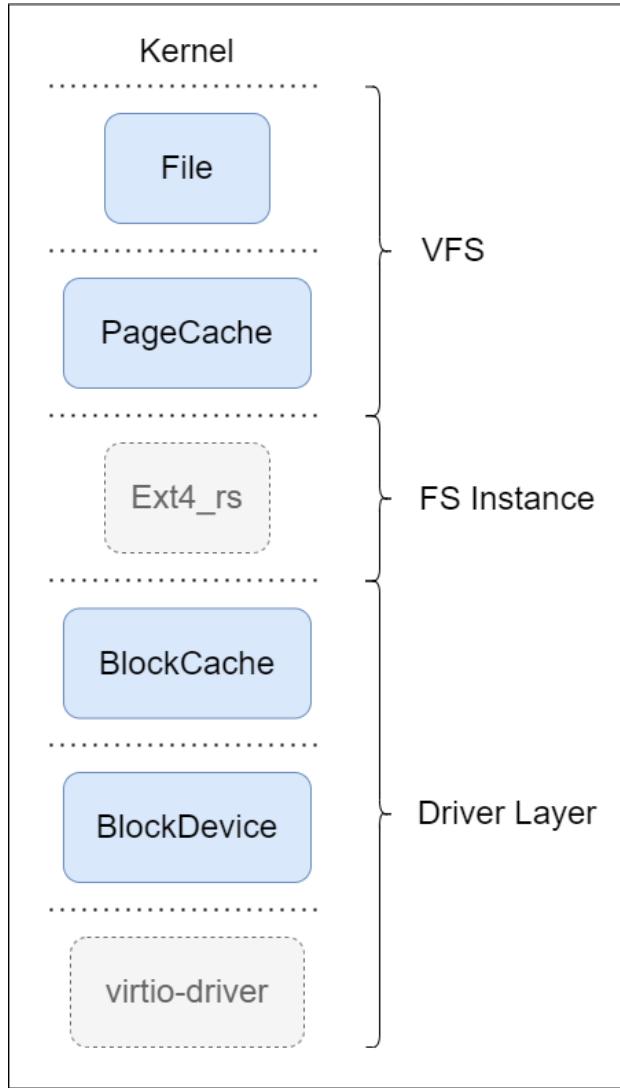


图 17: 缓存结构图

### 7.2.1 块缓存

在 NoAxiom 中，我们使用 `LruCache` 实现了异步块缓存 `BlockCache`，块缓存是驱动的缓存，用于加速驱动设备的读写。

```
1 pub struct AsyncBlockCache {
2     cache: SyncUnsafeCell<LruCache<usize, Arc<[u8; BLOCK_SIZE]>>>,
3     block_device: Arc<&'static dyn BlockDevice>,
4 }
5
6 #[async_trait]
7 impl BlockDevice for AsyncBlockCache {
8     fn device_name(&self) -> &'static str {
```

```

9     "AsyncBlockCache"
10    }
11    async fn read(&self, id: usize, buf: &mut [u8]) -> DevResult<usize> {
12        assert_eq!(buf.len(), BLOCK_SIZE);
13        let data = self.read_sector(id).await;
14        buf.copy_from_slice(&*data);
15        Ok(buf.len())
16    }
17    async fn write(&self, id: usize, buf: &[u8]) -> DevResult<usize> {
18        assert_eq!(buf.len(), BLOCK_SIZE);
19        let mut data: [u8; BLOCK_SIZE] = [0; BLOCK_SIZE];
20        data.copy_from_slice(buf);
21        self.write_sector(id, &data).await;
22        Ok(buf.len())
23    }
24    async fn sync_all(&self) -> DevResult<()> {
25        self.sync_all().await;
26        Ok(())
27    }
28}

```

其原理便是驱动设备每读取一个块号，都先检查块缓存是否已经存有副本，若有就可以直接从内存中读写，无需经过下层设备 I/O 操作，写操作需要标记块为 `Dirty`；若发现没有，就可以从内存中读取填充到块缓存中，再进行读写操作。因为块缓存容量有限，若填充时发现需要踢走已有的缓存，就需要考虑写回，我们仅当脏数据时写回。以上高效的设计，为我们的局部性反复读写文件功能提供了良好的性能。

### 7.2.2 页缓存

页缓存和块缓存本质相同，但所在层次不一致。`PageCache` 缓存的对象是文件在内存中打开后读写的信息，它可以加速更上层的请求，例如各类读写类系统调用。对于这类读写请求，可直接在 VFS 的层面获得所需数据，而无需通过具体文件系统访问，极大的提高了读写性能。

```

1 // NoAxiom/kernel/src/fs/pagecache.rs
2
3 #[derive(PartialEq)]
4 pub enum PageState {
5     Modified,
6     Shared,
7     Deleted,
8 }
9
10 pub struct Page {
11     data: FrameTracker,
12     state: PageState,
13 }
14
15 impl Page {
16     pub fn new(state: PageState) -> Self {
17         Self {
18             data: frame_alloc().unwrap(),
19             state,

```

```

20     }
21 }
22 pub fn as_mut_bytes_array(&mut self) -> &'static mut [u8] {
23     self.data.ppn().get_bytes_array()
24 }
25 pub fn mark_dirty(&mut self) {
26     if self.state != PageState::Deleted {
27         self.state = PageState::Modified;
28     }
29 }
30 pub fn mark_deleted(&mut self) {
31     self.state = PageState::Deleted;
32 }
33 }

```

我们使用 `Page` 作为页缓存中的基本存储单位。我们通过页帧管理器 `FRAME_ALLOCATOR` 申请页大小(4096B)的 `FrameTracker`，并且维护页缓存中每一个 `Page` 的状态：在写操作时进行脏标记，在删除文件时进行删除标记，为后面的写回策略提供依据。

```

1 // NoAxiom/kernel/src/fs/pagecache.rs
2 pub struct PageCacheManager {
3     inner: HashMap<Arc<dyn File>, HashMap<usize, Page>>,
4     page_count: usize,
5 }
6
7 impl PageCacheManager {
8     fn new() -> Self {
9         Self {
10             inner: HashMap::new(),
11             page_count: 0,
12         }
13     }
14
15     fn clean(&mut self, thresold: usize) {
16         ...
17     }
18
19     pub fn alloc(&mut self, state: PageState) -> Page {
20         ...
21         Page::new(state)
22     }
23
24     pub fn get_page(&self, file: &Arc<dyn File>, offset_align: usize) -> Option<&Page> {
25         assert!(is_aligned(offset_align, PAGE_SIZE));
26         if let Some(page_cache) = self.inner.get(file) {
27             if let Some(page) = page_cache.get(&offset_align) {
28                 return Some(page);
29             }
30         }
31         None
32     }
33
34     pub fn get_page_mut(&mut self, file: &Arc<dyn File>, offset_align: usize) -> Option<&mut Page> {
35         assert!(is_aligned(offset_align, PAGE_SIZE));

```

```

36     if let Some(page_cache) = self.inner.get_mut(file) {
37         if let Some(page) = page_cache.get_mut(&offset_align) {
38             return Some(page);
39         }
40     }
41     None
42 }
43
44 pub fn fill_page(&mut self, file: &Arc<dyn File>, offset_align: usize, page: Page)
45 {
46     assert!(is_aligned(offset_align, PAGE_SIZE));
47     self.inner
48         .entry(file.clone())
49         .or_insert_with(HashMap::new)
50         .insert(offset_align, page);
51 }
52
53 pub fn mark_deleted(&mut self, file: &Arc<dyn File>) {
54     self.inner.get_mut(file).map(|cache| {
55         cache.iter_mut().for_each(|(_, page)| {
56             page.mark_deleted();
57         })
58     });
59 }
60
61 pub fn truncate(&mut self, file: &Arc<dyn File>, length: usize) {
62     if let Some(page_cache) = self.inner.get_mut(file) {
63         let init_size = page_cache.len();
64         page_cache.retain(|&offset, _| offset < length);
65         let removed = init_size - page_cache.len();
66         self.page_count -= removed;
67     }
68 }
69
70 pub fn sync_all(&mut self) {
71     ...
72 }

```

以上是我们的页缓存管理结构的核心。`PageCacheManager` 存储着每个文件对应的已缓存内容，为 VFS 层面提供了读写接口，以及填充页缓存的回调函数。除此之外，为解决各类文件缓存一致性问题，我们还额外提供了 `truncate` 和 `mark_deleted` 接口，在需要调整文件大小和删除文件时使用。`sys_ftruncate` 和 `unlinkat` 系统调用便会用到这两个方法。

```

1 // NoAxiom/kernel/src/fs/pagecache.rs
2
3 lazy_static! {
4     pub static ref PAGE_CACHE_MANAGER: RwLock<PageCacheManager> =
5         RwLock::new(PageCacheManager::new());
6 }
7
8 #[inline(always)]
9 pub fn get_pagecache_rguard() -> RwLockReadGuard<'static, PageCacheManager> {
10     PAGE_CACHE_MANAGER.read()
11 }

```

```
12  
13 #[inline(always)]  
14 pub fn get_pagecache_wguard() -> RwLockWriteGuard<'static, PageCacheManager> {  
15     PAGE_CACHE_MANAGER.write()  
16 }
```

为尽可能提高高并发场景下对文件大量读取操作的性能，我们采用了读写锁 `RwLock` 进行对全局页缓存管理器的访问。

相较于往届队伍以及本届队伍，我们发现常规的设计中均为设计全局的缓存管理器，而是在 `Inode` 层面为每一个文件单独设立页缓存结构，然后在文件的生命周期彻底结束时（`Arc` 指针的引用计数为 0），进行脏数据对于原文件的写回并释放页缓存。

初赛中我们正是采用这样的设计，但在后续的研究中发现这样的设计只具有**局部视角**，因为当用户程序负载较大时，可能会有大量存在于内存中的文件描述符，并且每一个文件都具有页缓存，但所有文件在生命周期结束前都不会主动释放，这便会导致巨大的**内存压力**。如若此时又有更多的读写请求，极有可能导致页缓存申请失败，从而只能使用具体文件系统的读写方法。这不仅会导致文件系统层面的读写性能下降，更会影响到所有进程的执行速度（内存碎片过多，导致内存分配压力变大），严重时甚至会直接导致内核 panic！

基于以上分析与我们亲身实践经验，我们决定采用全局页缓存的设计进行整体缓存管理，并支持了缓存回收机制，保证内存的占用率水平。

```
// NoAxiom/kernel/src/fs/pagecache.rs

fn get_page_cache_capacity() -> usize {
    ...
    PAGE_CACHE_CAPACITY
}

fn get_page_cache_capacity_clean_threshold(page_cache_capacity: usize) -> usize {
    page_cache_capacity / 2
}

impl PageCacheManager {
    fn clean(&mut self, thresold: usize) {
        warn!("[PageCacheManager: clean], thresold: {}", thresold);
        let mut count = 0;
        for (file, page_cache) in self.inner.iter_mut() {
            debug!(
                "[PageCacheManager: clean] file: {}, page_cache size: {}",
                file.name(),
                page_cache.len()
            );
        }
        for (file, page_cache) in self.inner.iter_mut() {
            let file_size = file.size();
            let mut sorted = Vec::new();
            for (offset, page) in page_cache.iter() {
                sorted.push((offset, page));
            }
            sorted.sort_by_key(|(offset, _)| *offset); // for ext4_rs which has bug
            let mut offsets_to_remove = Vec::new();
```

```

32     for (offset, page) in &sorted {
33         if page.state == PageState::Modified {
34             assert_no_lock!();
35             assert!(Arch::is_external_interrupt_enabled());
36             let len = PAGE_SIZE.min(file_size - offset);
37             block_on(file.base_write(*offset, &page.data.ppn().get_bytes_array()
38                         ()[..len]))
39                         .unwrap();
40         }
41         offsets_to_remove.push(*offset);
42         count += 1;
43         if count >= thresold {
44             break;
45         }
46         for offset in offsets_to_remove {
47             page_cache.remove(&offset);
48         }
49         if count >= thresold {
50             self.page_count -= count;
51             return;
52         }
53     }
54 }
55
56 pub fn alloc(&mut self, state: PageState) -> Page {
57     let sys_capacity = get_page_cache_capacity();
58     let sys_thresold = get_page_cache_capacity_clean_threshold(sys_capacity);
59     if self.page_count > sys_capacity {
60         self.clean(sys_thresold);
61     }
62     self.page_count += 1;
63     Page::new(state)
64 }
65
66 pub fn sync_all(&mut self) {
67     for (file, page_cache) in self.inner.iter_mut() {
68         let file_size = file.size();
69         let mut sorted = Vec::new();
70         for (offset, page) in page_cache.iter() {
71             sorted.push((offset, page));
72         }
73         sorted.sort_by_key(|(offset, _)| *offset);
74         for (offset, page) in &sorted {
75             if page.state == PageState::Modified {
76                 assert_no_lock!();
77                 assert!(Arch::is_external_interrupt_enabled());
78                 let len = PAGE_SIZE.min(file_size - offset);
79                 match block_on(
80                     file.base_write(*offset, &page.data.ppn().get_bytes_array()[..
81                                     len]),
82                 ) {
83                     Ok(_) => {}
84                     Err(e) => {
85                         error!(
86                             "[PageCacheManager: sync_all] file: {}, offset: {},",

```

```

86             error: "{}",
87             file.name(),
88             offset,
89             e
90         );
91     }
92 }
93 page_cache.clear();
94 }
95 self.page_count = 0;
96 }
97 }

```

在每一次用户程序的申请页面请求向全局页缓存管理器发出时，页缓存管理器会首先通过 `get_page_cache_capacity()` 方法知道当前结点内存空间情况，以及清理阈值。如若发现当前页缓存使用容量 (`self.page_count`) 超过了动态可用阈值，便会执行 `clean()` 方法清理页缓存，然后再返回所需的页帧(此时的申请一定有效，并保证在内存给页缓存的可用空间内)。

总结来说 NoAxiom 的页缓存相较于传统方案具有以下几点优势：

- 低开销:** 因为我们采用了读写锁的设计，对于高频出现的读写操作，尤其是读操作，可以保证较高的 I/O 性能。
- 低内存占用:** 页缓存管理器会动态监测内存空间请求，保证全局页缓存的内存占用维持在动态阈值之内。我们在实践中还发现，良好的内存余量，有助于提高整体运行速度。
- 高性能:** 基于局部性原理，在可用容量内，我们会优先保证最新的文件的页缓存，保证整体的页缓存命中率维持在较高水平；实践中会频繁发现用户对于一个或某些特定文件的频繁开关、读写操作，传统设计中每一次文件生命周期结束时都会进行脏数据写回，也就对应着每一个 `sys_close` 系统调用都会执行一次文件写回，我们认为这是**低效的**。我们的全局页缓存设计不会对生命周期结束的文件进行立刻的写回，而是仍然保留，这样在下次打开时，便可以享受到之前已经读写的内容，可以省去大量的**具体文件系统操作**。而真正的写回只会发生在页缓存管理清理时。
- 高效一致性管理:** 我们为每一页标记了状态，对于处在 `Delete` 或者 `Share` 状态的页而言，可以直接丢弃，不需要执行具体文件系统的写回操作，我们仅对于处在 `Modified` 的页进行写回，最大程度的减少具体文件系统操作。

以下以读操作为例，展示了 `PageCache` 工作的原理：我们会将读请求按照页大小划分，每页的读取都会通过 rust 中 `copy_nonoverlapping` 的 `unsafe` 方法进行高效的复制；若出现页缺失现象，则会先调用 `base_read` 方法去文件系统中读出，再从页缓存中读出。

```

1 // NoAxiom/kernel/src/fs/vfs/basic/file.rs
2
3 pub async fn read_at(self: &Arc<dyn File>, offset: usize, buf: &mut [u8]) ->
4     SyscallResult {
5     let page_cache = self.page_cache();
6     if page_cache.is_none() {

```

```

6     assert_no_lock!();
7     return self.base_read(offset, buf).await;
8 }
9 let size = self.size();
10 if offset > size {
11     warn!(
12         "[read_at] {} offset {} > size {}",
13         self.meta().dentry.name()
14     );
15     return Ok(0);
16 }
17
18 let mut current_offset = offset;
19 let mut buf = buf;
20
21 // if the last page is not full,
22 // if the buf is run out
23 loop {
24     let (offset_align, offset_in) = align_offset(current_offset, PAGE_SIZE);
25     // maybe it is the last page, the exact len is equal to `size` -
26     // `current_offset` or the rest of the page
27     let page_len = (PAGE_SIZE - offset_in).min(size - current_offset);
28     // maybe the buf is not enough
29     let len = buf.len().min(page_len);
30
31     let r_guard = get_pagecache_rguard();
32     if let Some(page) = r_guard.get_page(&self.clone(), offset_align) {
33         unsafe {
34             core::ptr::copy_nonoverlapping(
35                 page.as_mut_bytes_array().as_ptr().add(offset_in),
36                 buf.as_mut_ptr(),
37                 len,
38             );
39         }
40
41         buf = &mut buf[len..];
42         current_offset += len;
43
44         if buf.is_empty() || current_offset == size {
45             break;
46         }
47
48         continue;
49     }
50     drop(r_guard);
51     let page_state = if self.inode().state() == InodeState::Deleted {
52         PageState::Deleted
53     } else {
54         PageState::Shared
55     };
56     let mut w_guard = get_pagecache_wguard();
57     let page = w_guard.alloc(page_state);
58     drop(w_guard);
59     self.base_read(offset_align, page.as_mut_bytes_array())
60         .await?;
61     let mut w_guard = get_pagecache_wguard();

```

```

62         w_guard.fill_page(&self.clone(), offset_align, page);
63     }
64
65     Ok((current_offset - offset) as isize)
66 }

```

### 7.2.3 路径缓存

我们在全局上添加了一个由自旋锁保护的路径缓存结构。其内部实现为 `HashMap`，记录了路径字符串到 `Dentry` 结构的映射，主要用于路径解析时的快速查找。

```

1 lazy_static::lazy_static! {
2     pub static ref PATH_CACHE: SpinLock<HashMap<String, Path>> = SpinLock::new(HashMap
3         ::new());
}

```

该路径缓存中保存的缓存生命周期为 `static`，即一旦有访问，便一直存在，直到用户通过 `unlinkat` 删除文件，则会从路径缓存中抹除。这样可以保证在复杂高负载场景下的路径缓存命中率，提高路径解析效率，并保证在 `open-close-open-close...` 序列中仍拥有良好的性能，而不是被 `sys_close` 系统调用去除缓存。

### 7.2.4 性能对比

	sync-driver	Interruptable-async
No-Cache	351	1026
BlockCache	1221	5463
BlockCache+PageCache	96210	99610

\*在高并发场景下的io测试

图 18: 不同条件下高并发场景的性能对比

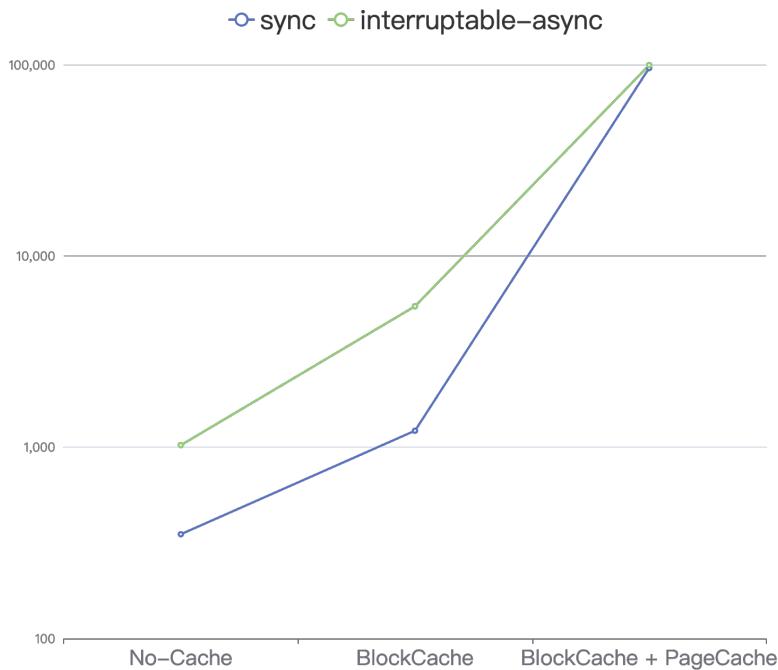


图 19: 不同条件下高并发场景的性能对比-折线图

以上数字均代表在高并发测试场景下的相对性能，**越大越好**。我们会发现在没有页缓存的高并发场景下，由异步驱动带来的**高并发性**使得整体性能优于传统驱动。但在添加了页缓存之后，由于大部分文件已经在内存中得到缓存，操作系统会通过页缓存直接读取，而不需要经过具体文件系统，所以下层的驱动带来的影响就不大了，但整体性能比没有页缓存的实现了**质的飞跃**。

	sync-driver	interruptable-async
No-Cache	1957	2851
BlockCache	311	296

\*单进程直接读写磁盘测试

图 20: 单进程磁盘读写性能

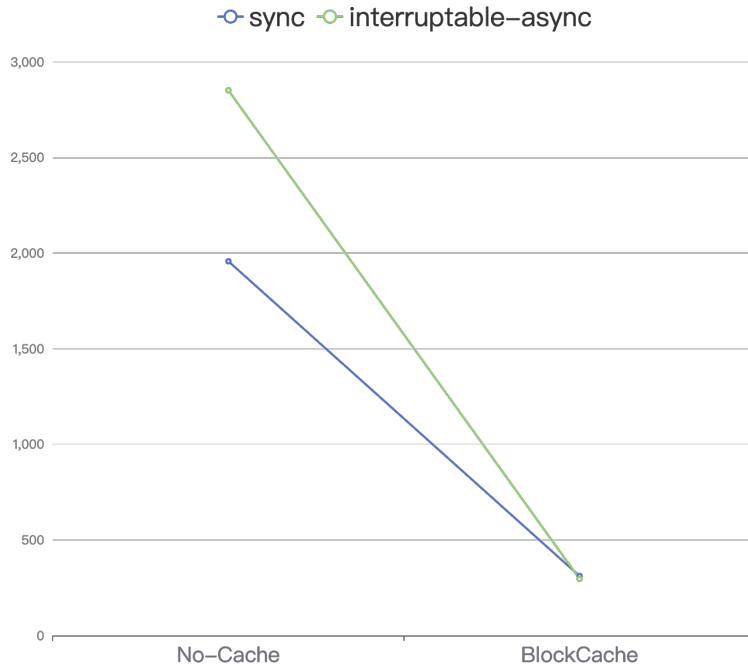


图 21: 单进程磁盘读写性能-折线图

我们还测试了一般场景下的读写性能，上面两个图描述了反复读写一片连续扇区时的时间消耗，**越小越好**。我们会发现在这种普通的单进程场景下，传统驱动表现反而比 NoAxiom 的异步驱动更好，原因就在于我们的异步驱动在单进程场景下仍然会让权，但此时调度器没有别的任务可以调度，只能继续调度本进程，这样反而造成了**时间粒度更大的忙等**，这里的时间粒度指的是异步调度器执行一次完整调度过程的开销，在这期间是无法立即响应的。而传统驱动通过简单的 `loop` 反而可以实现更快的响应。

但在加上块缓存优化后，由于是对于连续磁盘块的性能测试，此时块缓存就会缓存这些磁盘块，不管底层驱动是何种实现，大部分情况下都可以直接从块缓存中得到数据，也就导致了两种驱动的性能差距不大。

	sync-driver	interruptable-async
No-Cache	2245	3444
BlockCache	344	346
BlockCache+PageCache	132	123

\*单进程场景下的io测试

图 22: 单进程反复读写场景

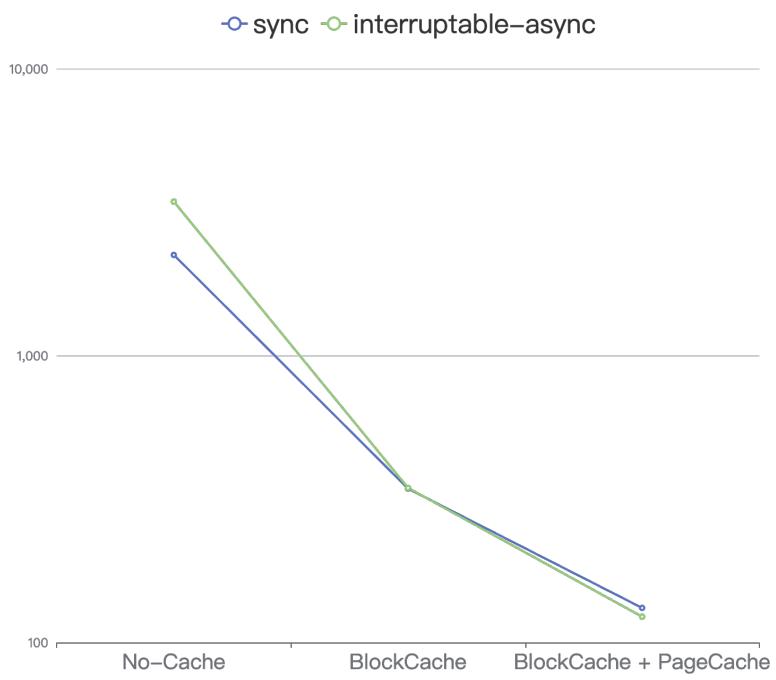


图 23: 单进程反复读写场景-折线图

以上两个图，更加综合的展示了块缓存和页缓存的适用场景和功能。它描述了单进程环境下反复读写文件的时间消耗，**越小越好**。在没有缓存的情况下，传统驱动反而比异步驱动优秀；在有块缓存或者页缓存的情况下，两者区别就不大了。

```

问题 124 输出 调试控制台 终端 端口 GITLENS GITLENS
Processor cache size set to 1024 kbytes.
Processor cache line size set to 32 bytes.
File stride size set to 17 * record size.
Throughput test with 4 processes
Each process writes a 1024 kbyte file in 1 kbyte records

Children see throughput for 4 initial writers = 35647.90 kB/sec
Parent sees throughput for 4 initial writers = 29506.73 kB/sec
Min throughput per process = 8306.16 kB/sec
Max throughput per process = 9623.67 kB/sec
Avg throughput per process = 8911.97 kB/sec
Min xfer = 960.00 kB

[PageCache] Hit 89887 / Total: 91984, HitRate: 96.85%

Children see throughput for 4 rewriters = 120825.96 kB/sec
Parent sees throughput for 4 rewriters = 61214.73 kB/sec
Min throughput per process = 0.00 kB/sec
Max throughput per process = 120825.96 kB/sec
Avg throughput per process = 30206.49 kB/sec
Min xfer = 0.00 kB

[PageCache] Hit 90111 / Total: 93008, HitRate: 96.89%

Children see throughput for 4 readers = 155269.14 kB/sec
Parent sees throughput for 4 readers = 83742.23 kB/sec
Min throughput per process = 0.00 kB/sec
Max throughput per process = 155269.14 kB/sec
Avg throughput per process = 38817.29 kB/sec
Min xfer = 0.00 kB

[PageCache] Hit 91135 / Total: 94032, HitRate: 96.92%

Children see throughput for 4 re-readers = 156383.62 kB/sec
Parent sees throughput for 4 re-readers = 82255.60 kB/sec
Min throughput per process = 0.00 kB/sec
Max throughput per process = 156383.62 kB/sec
Avg throughput per process = 39095.91 kB/sec
Min xfer = 0.00 kB

[PageCache] Hit 92159 / Total: 95056, HitRate: 96.95%

```

图 24: 页缓存命中率

```

386 RUN LTP CASE chmod01
387 [PageCache] Hit 39715 / Total: 43897, HitRate: 90.47%
388 tst_buffers.c:57: esc[1;34mTINFO: esc[0mTest is using guarded buffers
389 tst_memutils.c:152: esc[1;34mTINFO: esc[0moom_score_adj does not exist, skipping the adjustment
390 tst_test.c::1733: esc[1;34mTINFO: esc[0mLTP version: 20240524
391 tst_test.c::1617: esc[1;34mTINFO: esc[0mtimeout per run is 0h 00m 30s
392 tst_memutils.c:152: esc[1;34mTINFO: esc[0moom_score_adj does not exist, skipping the adjustment
393 chmod01.c:60: esc[1;34mTINFO: esc[0mTesting variant: verify permissions of file
394 chmod01.c:40: esc[1;32mTPASS: esc[0mchmod(testfile, 0000) passed
395 chmod01.c:50: esc[1;32mTPASS: esc[0mstat(testfile) mode=0000
396 chmod01.c:40: esc[1;32mTPASS: esc[0mchmod(testfile, 0007) passed
397 chmod01.c:50: esc[1;32mTPASS: esc[0mstat(testfile) mode=0007
398 chmod01.c:40: esc[1;32mTPASS: esc[0mchmod(testfile, 0070) passed
399 chmod01.c:50: esc[1;32mTPASS: esc[0mstat(testfile) mode=0070
400 chmod01.c:40: esc[1;32mTPASS: esc[0mchmod(testfile, 0700) passed
401 chmod01.c:50: esc[1;32mTPASS: esc[0mstat(testfile) mode=0700

```

图 25: 页缓存命中率

在官方测例 `iozone` 对于 I/O 的基准测试中，NoAxiom 的页缓存设计可以达到 **95%** 左右的优秀页缓存命中率；即使在官方测例 `ltp` 中，文件的读取纷繁复杂，命中率也可以达到 **90%** 左右。

## 7.3 文件系统实例

### 7.3.1 具体文件系统 EXT4/FAT32

我们为两种文件系统格式实现了上述 trait，其中 Ext4 文件系统，区别于往届传统队伍的 lwext4 库实现，我们使用了自主 fork 的 ext4\_rs<sup>13</sup>库，这是一个基于 rust 实现的异步文件系统<sup>14</sup>，提供了原库中所有支持操作的异步版本。

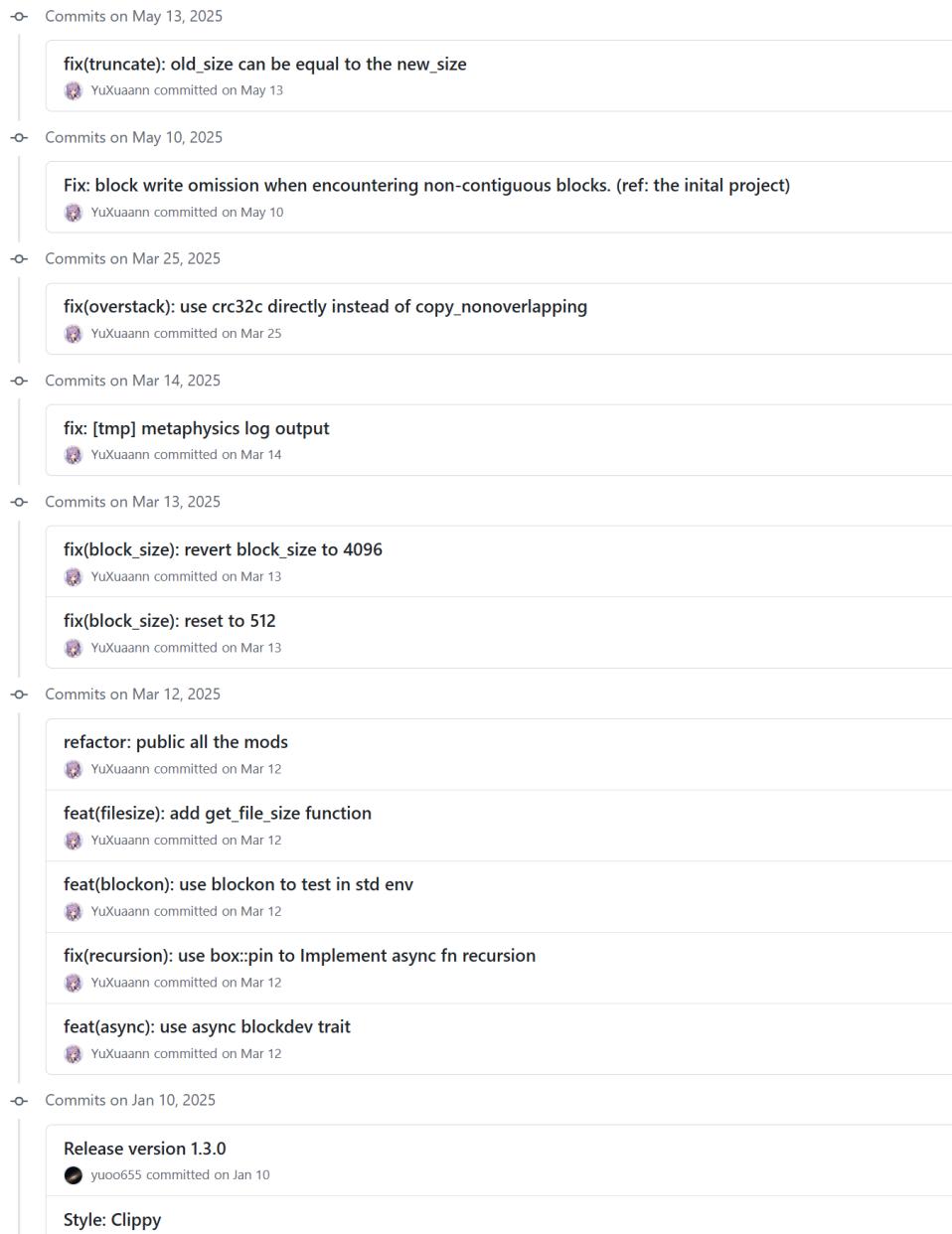


图 26: ext4\_rs-async 库的适配记录

<sup>13</sup>[https://github.com/yuoo655/ext4\\_rs](https://github.com/yuoo655/ext4_rs)

<sup>14</sup>[https://github.com/YuXuaann/ext4\\_rs-async-smp](https://github.com/YuXuaann/ext4_rs-async-smp)

### 7.3.2 管道文件

管道为进程（尤其是父子进程）之间的通信提供了良好的方法。其本质是一个生产者-消费者问题，我们需要解决多消费者和多生产者的互斥访问。

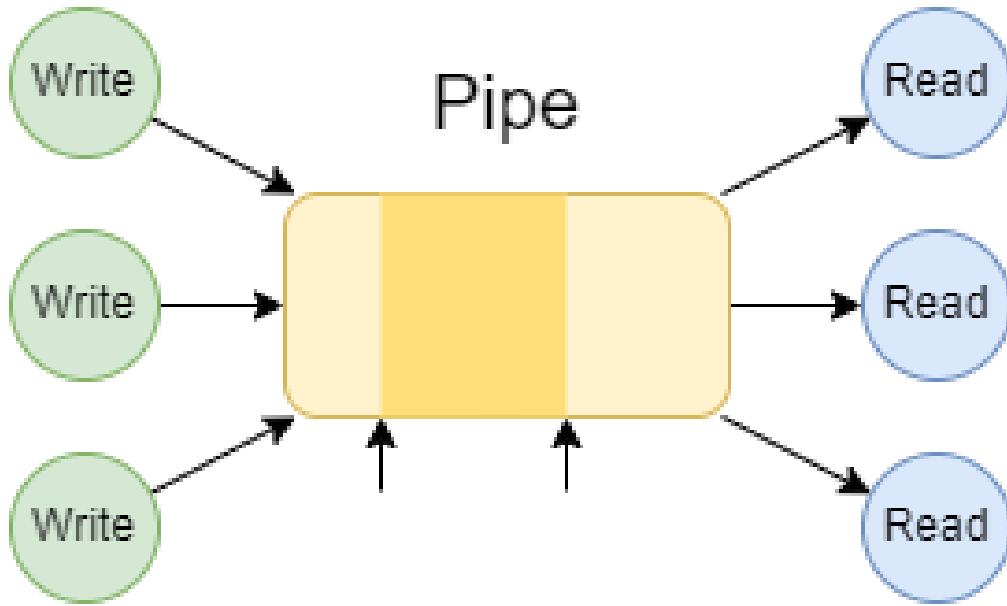


图 27: 管道

在传统内核的实现中，对于管道的访问采取“有多少读多少”，“能写多少写多少”的策略，进程之间不存在联系，只考虑当前管道中的内容，对于读请求如果管道为空就直接返回 0，对于写请求也是如果管道已满就返回零。因为用户程序的行为往往是通过循环等到足够的内容或空间进行读写，所以传统内核的做法导致了非常多空效果的系统调用，并且产生了大量切换特权态造成的上下文开销。

针对以上传统内核的弊端，NoAxiom 基于无栈协程机制设计了让权-唤醒机制的管道虚拟文件系统。我们为其实现了 VFS 的所有 trait，并且在创建新文件时，会虚拟挂载在根目录下，为其 inode 标记 FIFO。

```
1 pub struct PipeDentry {
2     meta: DentryMeta,
3 }
4
5 impl PipeDentry {
6     /// we mount all the pipes to the root dentry
7     pub fn new(name: &str) -> Arc<Self> {
8         let parent = root_dentry();
9         let super_block = parent.super_block();
10        let pipe_dentry = Arc::new(Self {
11            meta: DentryMeta::new(Some(parent.clone()), name, super_block),
12        });
13        debug!("[PipeDentry] create pipe dentry: {}", pipe_dentry.name());
14        parent.add_child_directly(pipe_dentry.clone());
15        pipe_dentry
16    }
}
```

17 }

对于一个管道读写请求，我们首先会检查管道 Buffer 中是否存在可供读写的内容，若没有，则将当前进程的 Waker 存到 PipeBuffer 结构体中，然后在 Future 中返回 Pending 让权。

```

1 struct PipeReadFuture {
2     read_len: usize,
3     pipe_buffer: Arc<SpinLock<PipeBuffer>>,
4 }
5 impl Future for PipeReadFuture {
6     // just for error handling
7     type Output = SysResult<()>;
8
9     fn poll(self: Pin<&mut Self>, cx: &mut Context<'_>) -> Poll<Self::Output> {
10         let mut buffer = self.pipe_buffer.lock();
11         if buffer.write_end {
12             if buffer.read_available() {
13                 Poll::Ready(Ok(()))
14             } else {
15                 buffer.add_read_event(self.read_len, cx.waker().clone());
16                 Poll::Pending
17             }
18         } else {
19             Poll::Ready(Ok(()))
20         }
21     }
22 }
```

若管道 Buffer 中存在可供读写的内容，首先进行管道的读写操作，然后唤醒存在 PipeBuffer 中的 Waker。这里在往届队伍中存在一种简单设计，即唤醒全部 Waker。我们在这里创新地设计了按需唤醒机制，我们只唤醒当前 PipeBuffer 中可用长度下能够唤醒的所有 Waker。因为每一个管道读写请求均存在长度参数，我们在存 Waker 的过程中还会一起记录此次读写请求的长度参数，当一次管道读写操作完成时，就会遍历 Waker 队列，pop 出管道 Buffer 长度可以 cover 住的 Waker 唤醒。

```

1 impl PipeBuffer {
2     fn notify_read_waker(&mut self) {
3         let mut read_available = self.read_available_len();
4         if read_available == 0 {
5             return;
6         }
7         while let Some((len, waker)) = self.read_wakers.pop() {
8             if read_available >= len {
9                 read_available -= len;
10                waker.wake();
11            } else {
12                waker.wake();
13                break;
14            }
15        }
16    }
17 }
18
19 #[async_trait]
20 impl File for PipeFile {
```

```

21     async fn base_read(&self, _offset: usize, buf: &mut [u8]) -> SyscallResult {
22         PipeReadFuture::new(buf.len(), self.buffer.clone()).await?;
23         let mut buffer = self.buffer.lock();
24         let ret = buffer.read(buf);
25         if ret != 0 {
26             buffer.notify_write_waker();
27         }
28         Ok(ret as isize)
29     }
30 }

```

总而言之，NoAxiom 的基于无栈协程机制的让权-唤醒机制，保证了不存在读写长度为 0 的管道读写系统调用，大大减少了上下文切换开销和系统调用数量。

顺带一提，NoAxiom 中提供给管道使用的 Buffer，是独立实现的 `RingBuffer` 结构，并且创新地采用了物理页帧进行存储，极大地避免了堆栈空间的占用。

```

1 // Ring buffer that max size is PIPE_BUF_SIZE
2 struct PipeBuffer {
3     data: FrameTracker,
4     head: usize,
5     tail: usize,
6     status: PipeBufferStatus,
7     read_wakers: Vec<(usize, Waker)>,
8     write_wakers: Vec<(usize, Waker)>,
9     read_end: bool,
10    write_end: bool,
11 }

```

### 7.3.3 套接字文件

首先，我们提供了枚举类 `Sock`，用来分发不同 socket。

```

1 pub enum Sock {
2     Tcp(TcpSocket),
3     Udp(UdpSocket),
4     // Unix(UnixSocket),
5 }
6
7 impl Sock {
8     pub fn bind(&mut self, addr: SockAddr, fd: usize) -> SysResult<()> {
9         let endpoint = addr.get_endpoint();
10        match self {
11            Sock::Tcp(socket) => socket.bind(endpoint, fd),
12            Sock::Udp(socket) => socket.bind(endpoint, fd),
13        }
14    }
15    pub fn listen(&mut self, backlog: usize) -> SysResult<()> {
16        match self {
17            Sock::Tcp(socket) => socket.listen(backlog),
18            _ => Err(Errno::ENOSYS),
19        }
20    }
21 // 下同

```

```

22     pub async fn connect(&mut self, addr: SockAddr) -> SysResult<()>
23     pub async fn accept(&mut self) -> SysResult<(TcpSocket, IpEndpoint)>
24     pub fn setsockopt(&mut self, level: usize, optname: usize, optval: &[u8]) ->
25         SysResult<()>
26     pub fn meta(&self) -> &SocketMeta
27     pub fn local_endpoint(&self) -> Option<IpEndpoint>
28     pub fn peer_endpoint(&self) -> Option<IpEndpoint>
29     pub async fn read(&mut self, buf: &mut [u8]) -> (SysResult<usize>, Option<
30         IpEndpoint>)
31     pub async fn write(&mut self, buf: &[u8], remote: Option<IpEndpoint>) -> SysResult<
32         usize>
33     pub fn shutdown(&mut self, operation: ShutdownType) -> SysResult<()>
34 }
```

然后我们通过 `AsyncMutex` 异步锁，包裹 `Sock` 结构体，就构成了套接字文件结构体。我们也将其视作普通文件，为其实现一系列标准的接口。

```

1  pub struct SocketFile {
2     meta: FileMeta,
3     sock: AsyncMutex<Sock>,
4     type_: PosixSocketType,
5 }
6
7 unsafe impl Send for SocketFile {}
8 unsafe impl Sync for SocketFile {}
9
10 #[async_trait]
11 impl File for SocketFile {
12     fn meta(&self) -> &FileMeta {
13         &self.meta
14     }
15     async fn base_read(&self, _offset: usize, buf: &mut [u8]) -> SyscallResult {
16         let mut sock = self.socket().await;
17         let res;
18         match &mut *sock {
19             Sock::Tcp(socket) => {
20                 res = socket.read(buf).await.0?;
21             }
22             Sock::Udp(socket) => {
23                 res = socket.read(buf).await.0?;
24             }
25         }
26         Ok(res as isize)
27     }
28     async fn base_write(&self, _offset: usize, buf: &[u8]) -> SyscallResult {
29         let mut sock = self.socket().await;
30         let res;
31         match &mut *sock {
32             Sock::Tcp(socket) => {
33                 res = socket.write(buf, None).await?;
34             }
35             Sock::Udp(socket) => {
36                 // use its remote endpoint as the destination
37                 res = socket.write(buf, None).await?;
38             }
39         }
40 }
```

```

40     }
41
42     Ok(res as isize)
43 }
44 fn poll(&self, req: &PollEvent, waker: Waker) -> PollEvent {
45     let mut sock = block_on(self.socket());
46     let poll_res = match &mut *sock {
47         Sock::Tcp(socket) => socket.poll(req, waker),
48         Sock::Udp(socket) => socket.poll(req, waker),
49     };
50
51     let mut res = PollEvent::empty();
52     if poll_res.contains(PollEvent::POLLIN) {
53         res |= PollEvent::POLLIN;
54     }
55     if poll_res.contains(PollEvent::POLLOUT) {
56         res |= PollEvent::POLLOUT;
57     }
58     if poll_res.contains(PollEvent::POLLHUP) {
59         warn!("[Socket::poll] PollEvent is hangup");
60         res |= PollEvent::POLLHUP;
61     }
62     res
63 }
64 }

```

套接字文件和管道文件一样，也提供了虚拟挂载在根目录的功能。关于套接字文件具体的功能实现，在后文的网络模块中有详细阐述。

#### 7.3.4 基于内存的文件系统 RamFs

RamFs 是建立在内存上的非持久文件系统，使用堆内存创建的 `Vec` 去存储文件的内容。该文件系统目前仅用来实现部分虚拟文件，用于满足用户交互的需求。具有响应速度快的优势。

### 7.4 虚拟文件

凭借完善的 VFS 结构设计，下面的各类文件实例，只需要实现相应的 `trait` 即可。

#### 7.4.1 进程管理 /proc

`/proc` 用来存储内核当前的运行状态，用户可以通过 `proc` 目录下的文件查看当前正在运行进程的信息，也可以通过与这些文件交互实现和内核进程的直接交互。

1. `/proc/meminfo`：读取时返回内存信息。
2. `/proc/mounts`：读取时返回挂载的文件系统信息。
3. `/proc/sys/`：基于 RamFs，存储关于系统相关文件信息。
4. `/proc/interrupts`：全局中断统计文件。

5. `/proc/sys/kernel/pid_max` : 基于 RamFs, 读取时返回最大进程号。
6. `/proc/self/*` : 基于 RamFs, 存储关于当前进程相关信息。
7. `/proc/{x}/*` : 基于 RamFs, `x` 表示进程号, 存储关于该进程相关信息。
8. `/proc/{x}/exe` : `x` 代表 self 当前进程或者其他进程号, 用 `readlink` 读取时返回该进程正在执行的文件的绝对路径。
9. `/proc/{x}/maps` : `x` 代表 self 当前进程或者其他进程号, 描述了进程关于 `mmap` 区域的相关信息。
10. `/proc/{x}/stat` : `x` 代表 self 当前进程或者其他进程号, 描述了进程相关信息。
11. `/proc/{x}/status` : `x` 代表 self 当前进程或者其他进程号, 描述了进程相关信息。
12. `/proc/{x}/fd/{y}` : `x` 代表 self 当前进程或者其他进程号, `y` 表示某个文件描述符, `fd` 目录表示该进程的文件描述符表, 里面的内容为各个文件的文件描述符, 使用 `sys_read` 进行读取时, 等价于读取那个文件。

#### 7.4.2 虚拟设备 /dev

`/dev` 用来存储那些代表硬件设备的文件, 它允许用户使用与文件交互的方式与硬件交互。

1. `/dev/null` : 垃圾桶。读操作返回空, 写操作相当于丢弃数据。
2. `/dev/zero` : 读取时会返回全 0.
3. `/dev/tty` : 非常重要的字符设备。它代表当前正在使用的控制终端设备, 目前也是我们内核的串口设备, 这意味着实际上每个进程初始化时 `FdTable` 中默认的 `STDIN`、`STDOUT`、`STDERR` 的实际实现都是 `TtyFile` ; 我们也为其实现了全面的 `ioctl` 方法。
4. `/dev/rtc` : 读取时返回实时时钟。
5. `/dev/cpu_dma_latency` : 读取时返回 CPU 的 DMA 延迟。
6. `/dev/urandom` : 读取时返回随机数。
7. `/dev/shm` : 用于实现进程间的共享内存。
8. `/dev/loop-control` : 回环设备相关虚拟文件, 用于用户以一个普通文件虚拟一个块设备。

## 7.5 系统调用实现

```
// fs
SYS_FCHMOD => self.sys_fchmod(fd: args[0], mode: args[1]),
SYS_FCHMODAT => self.sys_fchmodat(fd: args[0], path: args[1], mode:
SYS_FCHOWN => self.sys_fchown(fd: args[0], owner: args[1] as u32,
SYS_FCHOWNAT => self.sys_fchownat(fd: args[0], path: args[1], owner:
SYS_READ => self.sys_read(fd: args[0], buf: args[1], len: args[2]),
SYS_READV => self.sys_readv(fd: args[0], iovp: args[1], iovcnt: args[2]),
SYS_PREAD64 => self.sys_pread64(fd: args[0], buf: args[1], len: args[2]),
SYS_PREADV => self.sys_preadv(fd: args[0], iovp: args[1], iovcnt: args[2]),
SYS_PREADV2 => self.sys_preadv2(fd: args[0], iovp: args[1], iovcnt: args[2]),
SYS_WRITE => self.sys_write(fd: args[0], buf: args[1], len: args[2]),
SYS_WRITEV => self.sys_writev(fd: args[0], iovp: args[1], iovcnt: args[2]),
SYS_PWRITE64 => self.sys_pwrite64(fd: args[0], buf: args[1], len: args[2]),
SYS_PWRITEV => self.sys_pwritev(fd: args[0], iovp: args[1], iovcnt: args[2]),
SYS_PWRITEV2 => self.sys_pwritev2(fd: args[0], iovp: args[1], iovcnt: args[2]),
SYS_CLOSE => self.sys_close(fd: args[0]),
SYS_MKDIRAT => self.sys_mkdirat(dirfd: args[0] as isize, path: args[1]),
SYS_OPENAT => self.sys_openat(fd: args[0] as isize, path: args[1]),
SYS_CHDIR => self.sys_chdir(path: args[0]),
SYS_FCHDIR => self.sys_fchdir(fd: args[0]),
SYS_CHROOT => self.sys_chroot(path: args[0]),
SYS_GETCWD => self.sys_getcwd(buf: args[0], size: args[1]).await,
SYS_DUP => self.sys_dup(fd: args[0]),
SYS_DUP3 => self.sys_dup3(old_fd: args[0], new_fd: args[1], flags: args[2]),
SYS_PIPE2 => self.sys_pipe2(pipe: args[0], flags: args[1] as i32),
SYS_STATX => self.sys_statx(dirfd: args[0] as isize, path: args[1],
SYS_FSTAT => self.sys_fstat(fd: args[0], stat_buf: args[1]).await,
SYS_FSTATAT => self.sys_fstatat(dirfd: args[0] as isize, path: args[1],
SYS_STATFS => self.sys_statfs(path: args[0], buf: args[1]).await,
SYS_FSTATFS => self.sys_fstatfs(fd: args[0], buf: args[1]).await,
SYS_GETDENTS64 => self.sys_getdents64(fd: args[0], buf: args[1], len: args[2]),
SYS_MOUNT => self.sys_mount(special: args[0], dir: args[1], fstype: args[2]),
SYS_UMOUNT2 => self.sys_umount2(dir: args[0], flags: args[1]),
SYS_LINKAT => self.sys_linkat(olddirfd: args[0] as isize, oldpath: args[1],
SYS_SYMLINKAT => self.sys_symlinkat(target: args[0], newdirfd: args[1],
SYS_UNLINKAT => self.sys_unlinkat(dirfd: args[0] as isize, path: args[1]),
SYS_PRLIMIT64 => self.sys_prlimit64(pid: args[0], resource: args[1] as i32, rlimit: args[2]),
SYS_GETRLIMIT => self.sys_getrlimit(resource: args[0] as _, rlimit: args[1]),
SYS_SETRLIMIT => self.sys_setrlimit(resource: args[0] as _, rlimit: args[1]),
SYS_FCNTL => self.sys_fcntl(fd: args[0], cmd: args[1], arg: args[2]),
SYS_READLINKAT => self.sys_readlinkat(dirfd: args[0] as isize, path: args[1]),
SYS_IOCTL => self.sys_ioctl(fd: args[0], request: args[1], arg: args[2]),
SYS_SENDFILE => self.sys_sendfile(out_fd: args[0], in_fd: args[1], offset: args[2]),
SYS_FACCESSAT => self.sys_faccessat(fd: args[0], path: args[1], mode: args[2]),
SYS_FACCESSAT2 => self.sys_faccessat2(fd: args[0], path: args[1], mode: args[2]),
SYS_UTIMENSAT => self.sys_utimensat(dirfd: args[0] as isize, path: args[1],
SYS_LSEEK => self.sys_lseek(fd: args[0], offset: args[1] as isize, whence: args[2]),
SYS_RENAMEAT2 => self.sys_renameat2(old_dirfd: args[0] as isize, oldpath: args[1],
SYS_COPY_FILE_RANGE => self.sys_copy_file_range(fd_in: args[0], off_in: args[1],
SYS_TRUNCATE => self.sys_truncate(path: args[0], len: args[1] as isize),
SYS_FTRUNCATE64 => self.sys_ftruncate(fd: args[0], length: args[1]).await,
SYS_SPLICE => self.sys_splice(fd_in: args[0], off_in: args[1], fd_out: args[2], off_out: args[3]),
SYS_FALLOCATE => self.sys_fallocate(fd: args[0], mode: args[1] as i32, offset: args[2], len: args[3]),
SYS_FADVISE64 => self.sys_fadvise64(fd: args[0], offset: args[1], len: args[2], flags: args[3]),
SYS_MKNODAT => self.sys_mknodat(fd: args[0] as isize, path: args[1], mode: args[2]),
SYS_SYNC_FILE_RANGE => self.sys_sync_file_range(fd: args[0], offset: args[1], len: args[2])
```

图 28: fs 相关系统调用实现

### 7.5.1 I/O 读写类

上面在缓存中已经介绍了大致的读写操作。但在实际读写系统调用中，我们还需要通过信号机制中提供的 `intable` 接口，为那些可能产生阻塞的读写操作（例如 `socket` 读写或管道读写）提供强制退出的可能。

### 7.5.2 I/O 多路复用

在 `File` 结构体中，我们强制要求每一个文件都需要实现 `poll` 方法，就是为了优雅的实现 I/O 的多路复用功能。通过 rust 强大的 `trait` 功能，我们只需要使用 VFS 层提供的抽象接口，便可实现对所有文件（可以是各种不同类型的文件）的多路监听。

其中 `poll` 方法原型是这样定义的：

```
1 fn poll(&self, req: &PollEvent, waker: Waker) -> PollEvent;
```

它会传入当前进程的 `Waker`，以及 `poll` 操作的具体请求，具体文件需要返回 `poll` 请求中它可以接受处理的那些请求。如果某些请求暂时还无法得到处理，具体文件就会将当前进程的 `Waker` 存入它内部的 `Buffer` 中。

对于具体系统调用来说，如果它观察到文件 list 中无一准备就绪，调度器就会执行让权操作，执行其他进程的任务，直到具体文件唤醒原进程。以上操作只需要实现 `Future trait` 即可实现，这样的思路完美利用了无栈协程机制，以高效优雅的方式实现了 I/O 多路复用功能。

```
1 pub struct PpollItem {
2     id: usize,
3     events: PollEvent,
4     file: Arc<dyn File>,
5 }
6
7 pub struct PpollFuture {
8     fds: Vec<PpollItem>,
9 }
10
11 impl Future for PpollFuture {
12     type Output = Vec<(usize, PollEvent)>;
13
14     fn poll(self: Pin<&mut Self>, cx: &mut Context<'_>) -> Poll<Self::Output> {
15         let mut result = Vec::new();
16         for poll_item in &self.fds {
17             let id = poll_item.id;
18             let req = &poll_item.events;
19             let res = poll_item.file.poll(req, cx.waker().clone());
20             if !res.is_empty() {
21                 result.push((id, res));
22             }
23         }
24         if result.is_empty() {
25             Poll::Pending
26         } else {
27             Poll::Ready(result)
28         }
29     }
30 }
```

```
29     }
30 }
```

特别的，我们对于系统调用 `sys_pselect` 也复用了上述 `Future` 实现。

## 7.6 其他相关结构

### 7.6.1 文件描述符表 FdTable

每一个进程都拥有一个文件描述符表。因为在用户的视角，文件是一个抽象的概念，用户是通过文件描述符 `fd` (一个非负整数) 来操控文件的，而内核中进程持有的文件描述符表就是维护这一映射关系的结构。`FdTable` 记录了每一个 `fd` 对应的文件及其 `FcntlArgFlags`。在进程 `fork` 时，文件描述符表会完整复制，我们只需要为 `FdTable` 标记 `##[derive(Clone)]` 就可以优雅地自动实现 `Fdtable` 复制方法。

```
1 #[derive(Clone)]
2 pub struct FdTableEntry {
3     flags: FcntlArgFlags,
4     pub file: Arc<dyn File>,
5 }
6
7 #[derive(Clone)]
8 pub struct FdTable {
9     pub table: Vec<Option<FdTableEntry>>,
10    rlimit: RLimit,
11 }
```

其中 `FcntlArgFlags` 结构主要是用来标记该文件描述符是否拥有 `close_on_exec` 标志位，该标志位可以保证这个文件描述符在执行了 `sys_execve` 后被关闭，防止出现安全问题。

`FdTable` 还承担了对于 `RLimit` 的管理：

```
1 /// Resource Limit from linux
2 #[derive(Debug, Clone, Copy)]
3 pub struct RLimit {
4     /// Soft limit
5     pub rlim_cur: usize,
6     /// Hard limit (ceiling for rlim_cur)
7     pub rlim_max: usize,
8 }
```

若 `alloc` 操作申请的 `fd` 超过了最大描述符数量的限制，便会返回 `Errno::EMFILE`。除此之外，也为系统调用 `sys_prlimit` 提供了设置当前进程文件描述符的 `RLlimit` 接口：

```
1 pub fn rlimit_mut(&mut self) -> &mut RLimit {
2     &mut self.rlimit
3 }
4
5 pub async fn sys_prlimit64(
6     &self,
7     pid: usize,
8     resource: u32,
9     new_limit: usize,
10    old_limit: usize,
```

```

11 ) -> SyscallResult {
12     let task = ...
13     let mut fd_table = task.fd_table();
14     let new_limit = UserPtr::<RLimit>::new(new_limit);
15     ...
16     if !new_limit.is_null() {
17         match resource {
18             Resource::NOFILE => *fd_table.rlimit_mut() = new_limit.read().await?,
19             _ => {}
20         }
21     }
22     Ok(0)
23 }
```

为充分利用内存资源和提高内核效率，若用户使用了 `sys_close` 关闭了某个文件描述符 `fd` 对应的文件，那么该 `fd` 对应的槽位将仅需要设置为 `None`，性能开销极小。下次申请 `fd` 时会优先选择空槽位，提高了内存的利用率。

```

1 // get the soft rlimit of the fd table
2 fn rslimit(&self) -> usize {
3     self.rlimt.rlim_cur
4 }
5
6 /// Allocate a new fd slot, if has empty slot, return the first empty slot,
7 /// you should use the fd after alloc immediately
8 pub fn alloc_fd(&mut self) -> SysResult<usize> {
9     if let Some(fd) = self.table.iter().position(|x| x.is_none()) {
10         return Ok(fd);
11     }
12
13     if self.table.len() >= self.rslimit() {
14         Err(Errno::EMFILE)
15     } else {
16         self.table.push(None);
17         Ok(self.table.len() - 1)
18     }
19 }
```

## 7.7 文件一致性管理

在 `Inode` 生命周期结束时，会将 `PageCache` 写回；在系统 `Shutdown` 时，显式调用 `sync` 方法，写回脏数据。

## 7.8 文件权限管理

文件权限是一个涉及 `Task` 模块和文件系统的综合设计，它需要综合考量访问行为、文件本体权限、用户身份权限这三大因素。相较于传统设计中在 `Task` 模块中直接添加若干 `AtomicUsize` 成员用于权限管理（例如 `uid`, `gid`, `fsuid...`），对于每一个 `id` 都使用一个原子变量，`NoAxiom` 使用了 `TaskUserId` 结构体进行统一管理：

### 7.8.1 权限结构设计

因为考虑到文件权限检查中，经常涉及这些权限中的若干值，我们使用统一结构体并且上锁读取，只需一次原子指令开销，而非原先若干原子指令开销。提高了高频发生的权限检查的性能。

```
1 // NoAxiom/kernel/src/task/taskid.rs
2 #[derive(Debug, Clone)]
3 pub struct TaskUserId {
4     /// user id
5     pub uid: u32,
6     /// group id
7     pub gid: u32,
8     /// user id - file system
9     pub fsuid: u32,
10    /// group id - file system
11    pub fsgid: u32,
12    /// user id - effective
13    pub euid: u32,
14    /// group id - effective
15    pub egid: u32,
16    /// user id - saved
17    pub suid: u32,
18    /// group id - saved
19    pub sgid: u32,
20 }
```

对于 Inode 而言，使用了原子变量压缩保存所有权限和文件类型位，使得读取开销最低。

```
1 // NoAxiom/kernel/src/fs/vfs/basic/inode.rs
2
3 pub struct InodeMeta {
4     ...
5     pub inode_mode: AtomicU32,
6 }
```

### 7.8.2 权限检查功能

具体功能实现上，内核统一使用 Dentry 提供的 `check_access()` 方法进行检查（部分参考 Rocket 队）：

```
1 // NoAxiom/kernel/src/fs/vfs/basic/dentry.rs
2
3 pub fn can_search(self: &Arc<Self>, task: &Arc<Task>) -> bool {
4     ...
5
6     let (user_perm, group_perm, other_perm) =
7         ((i_mode >> 6) & 0o7, (i_mode >> 3) & 0o7, i_mode & 0o7);
8     let perm = if euid == self.inode().unwrap().uid() {
9         user_perm
10    } else if egid == self.inode().unwrap().gid() {
11        group_perm
12    } else {
13        other_perm
14    };
15 }
```

```

16     if perm & 0o111 == 0 {
17         return false;
18     }
19     true
20 }
21
22 pub fn check_access(
23     self: &Arc<Self>,
24     task: &Arc<Task>,
25     access: i32,
26     is_fs: bool,
27 ) -> SysResult<()> {
28     let user_id = task.user_id();
29     let (uid, gid) = if is_fs {
30         (user_id.fsid(), user_id.fsgid())
31     } else {
32         (user_id.uid(), user_id.gid())
33     };
34
35     ...
36
37     if uid == UID_ROOT {
38         if (access & X_OK != 0) && (pri.bits() & 0o111 == 0) {
39             error!(
40                 "[sys_faccessat] root user cannot execute file: {:?}",
41                 self.path()
42             );
43             return Err(Errno::EACCES);
44         }
45         return Ok(());
46     }
47
48     // check if the parent directory is accessible
49     if let Some(parent) = dentry.parent() {
50         parent.check_arrive()?;
51     }
52
53     if access == F_OK {
54         return Ok(());
55     }
56
57     let permission = if uid == inode.uid() {
58         pri.user_permissions() as i32
59     } else if gid == inode.gid() {
60         pri.group_permissions() as i32
61     } else {
62         pri.other_permissions() as i32
63     };
64
65     if (access & X_OK != 0) && (permission & X_OK == 0) {
66         error!("[check_access] user cannot execute file: {:?}", self.path());
67         return Err(Errno::EACCES);
68     }
69
70     if (access & R_OK != 0) && (permission & R_OK == 0) {
71         error!("[check_access] user cannot read file: {:?}", self.path());

```

```
72     return Err(Errno::EACCES);
73 }
74
75 if (access & W_OK != 0) && (permission & W_OK == 0) {
76     error!("[check_access] user cannot write file: {:?}", self.path());
77     return Err(Errno::EACCES);
78 }
79 Ok(())
80 }
```

除此之外也提供了检查路径可达性的方法，通过完善的 Dentry 设计，检查路径上文件夹的可执行位是否设置：

```
1 // NoAxiom/kernel/src/fs/vfs/basic/dentry.rs
2
3 pub fn check_arrive(self: &Arc<Self>) -> SysResult<()> {
4     if self.inode()?.privilege().bits() & 0o111 == 0 {
5         return Err(Errno::EACCES);
6     }
7     if let Some(parent) = self.parent() {
8         return parent.check_arrive();
9     }
10    Ok(())
11 }
```

# 八、信号系统

## 8.1 信号管理

### ► 待处理信号管理

NoAxiom 使用 `PendingSigs` 维护等待处理的信号，其内部维护了当前进程对应的待处理信号信息。

```
1 pub struct SigPending {
2     pub sig_mask: SigMask,           // signal mask of the task
3     pub queue: VecDeque<SigInfo>, // pending signal queue that should be handled
4     pub pending_set: SigSet,        // current pending signal set
5     pub should_wake: SigSet,        // signals that should wake the task
6 }
```

我们使用 `SpinLock` 自旋锁保护 `PendingSigs`，使其访问具备原子性。

### ► 信号捕获函数管理

我们使用 `SigAction` 对于某一信号的具体处理逻辑进行管理。其中 `SAHandlerType` 字段给定了当前处理函数由用户定义还是执行默认操作，若由用户定义，则在检查信号时会进行用户信息的准备。

```
1 pub struct KSigAction {
2     pub handler: SAHandlerType,
3     pub mask: SigMask,
4     pub flags: SAFlags,
5     pub restorer: usize,
6 }
```

### ► 信号具体信息维护

为了便于维护信号信息，我们使用 `SigInfo` 对于信号传递的信息进行了详细描述。其中包括了：信号编号、信号类型、信号错误码、信号详细信息等信息。其中 `detail` 表示了信号详细信息，如 `SIG_KILL` 发生时，我们将通过该字段传递发送该信号的具体进程，以此便于内核检查。

```
1 pub struct SigInfo {
2     pub signo: Signo,
3     pub code: SigCode,
4     pub errno: SigErrno,
5     pub detail: SigDetail,
6 }
```

## 8.2 信号处理流程

### ► 信号检测

在任务运行的主函数 `task_main` 中，在执行完用户系统调用后，我们会对于当前进程待处理的信号进行检测。其具体处理逻辑为：检查是否具有等待处理的信号，并访问当前进程的 `SigAction` 进行具体的信号处理。假如当前的 `SigAction` 由用户定义，那么还需要进行对应的用户上下文与用户栈信息的更新。

```
1 task.check_signal(None).await;
```

## ► 信号返回

当需要执行 `sigreturn` 系统调用时，由于需要在用户态执行系统调用，因此需要对于对应跳板页进行用户空间的映射。我们在地址空间初始化的时候就对于该页面进行了映射。具体返回时，我们将根据不同的架构内嵌一段汇编代码指示用户进行返回。

```
1 .section .text.signal
2 .globl user_sigreturn
3 .align 12
4 user_sigreturn:
5 ori      $a7, $zero, 139          # syscall SIGRETURN
6 syscall 0
```

## 8.3 可中断系统调用

Linux 中规定了部分系统调用可被特定信号打断，如 `futex`、`ppoll` 等调用均会在调用过程中接受信号时，返回错误码 `Errno::EINTR`，以防进程因系统调用永久挂起而导致进程卡死。

为了统一规范该行为，NoAxiom 中设计了 `interruptable(future)` 函数用于进行可打断操作的异步状态机构建。具体而言，调用该函数时，会将原有的异步操作在外部进一步包装为 `IntableFuture`，该异步函数会在内部 `future` 返回 `Pending` 的时候，检测当前进程的信号状态，假如检测到存在等待处理的信号，则直接返回 `Errno::EINTR` 用于指示当前异步函数被打断。

```
1 impl<F, T> Future for IntableFuture<'_, F>
2 where
3     F: Future<Output = T>,
4 {
5     type Output = SysResult<T>;
6     fn poll(self: Pin<&mut Self>, cx: &mut Context<'_>) -> Poll<Self::Output> {
7         let this = self.project();
8         let task = this.task;
9         match this.fut.poll(cx) {
10             Poll::Ready(res) => Poll::Ready(Ok(res)),
11             Poll::Pending => {
12                 // start to handle signal
13                 let mask = this.mask;
14                 let pcb = task.pcb();
15                 if pcb.signals.has_pending_signals(*mask) {
16                     task.record_current_result_reg();
17                     task tcb_mut().flags |= TaskFlags::TIF_SIGPENDING;
18                     Poll::Ready(Err(Errno::EINTR))
19                 } else {
20                     Poll::Pending
21                 }
22             }
23         }
24     }
25 }
26 pub async fn interruptable<T>(
27     task: &Arc<Task>,
28     fut: impl Future<Output = T>,
29     new_mask: Option<SigMask>,
30     wake_set: Option<SigSet>,
31 ) -> SysResult<T> {
```

```

32     // set new mask
33     if let Some(new_mask) = new_mask {
34         task.swap_in_sigmask(new_mask.without_kill());
35         task.tif_mut().insert(TaskFlags::TIF_RESTORE_SIGMASK);
36     }
37     let mask = task.sig_mask();
38
39     // set wake signal: forbid masked and ignored signals
40     let wake_set = wake_set.unwrap_or(SigSet::empty());
41     let ignored_set = task.sa_list().get_ignored_bitmap();
42     let wake_set = !(mask | ignored_set) | wake_set;
43     task.pcb().set_wake_signal(wake_set);
44
45     // suspend now!
46     let res = IntableFuture { task, fut, mask }.await;
47
48     // restore old mask
49     task.pcb().clear_wake_signal();
50     res
51 }
```

例如，在 `ppoll` 实现中，我们使用 `interruptable` 包装了原有的轮询操作，使其变为可打断形式，并通过 `interruptable(...).await?` 实现了错误码的传回。该函数极大地简化了各类可中断系统调用的处理过程，实现了可中断调用的处理逻辑统一化，极大地增加了代码的简洁性。

```

1 // ppoll
2 let interruptable = interruptable(self.task, fut, None);
3 let res = match interruptable.await? {
4     TimeLimitedType::Ok(res) => res,
5     TimeLimitedType::TimeOut => return Ok(0),
6 };
```

此外，为了便于异步检测任务是否正在处于信号处理阶段，我们参照 Linux 中的实现，为 `Task` 添加了 `TaskFlags` 用于进行异步的任务状态通知。当任务存在待处理的信号时，会设置 `TIF_SIGPENDING` 位以通知任务进行信号处理。此外，对于需要恢复 `sigmask` 的系统调用，我们还添加了 `TIF_RESTORE_SIGMASK` 用于进行异步的掩码恢复过程。

```

1 bitflags! {
2     /// Thread flags defined to check the state of a thread asynchronously.
3     pub struct TaskFlags: usize {
4         /// restore sig mask in do_signal()
5         const TIF_RESTORE_SIGMASK = 1 << 0;
6         /// resumption notification requested
7         const TIF_NOTIFY_RESUME = 1 << 1;
8         /// signal pending
9         const TIF_SIGPENDING = 1 << 2;
10        /// rescheduling necessary
11        const TIF_NEED_RESCHED = 1 << 3;
12        /// signal notifications exist
13        const TIF_NOTIFY_SIGNAL = 1 << 5;
14    }
15 }
```

# 九、网络模块

## 9.1 结构概述

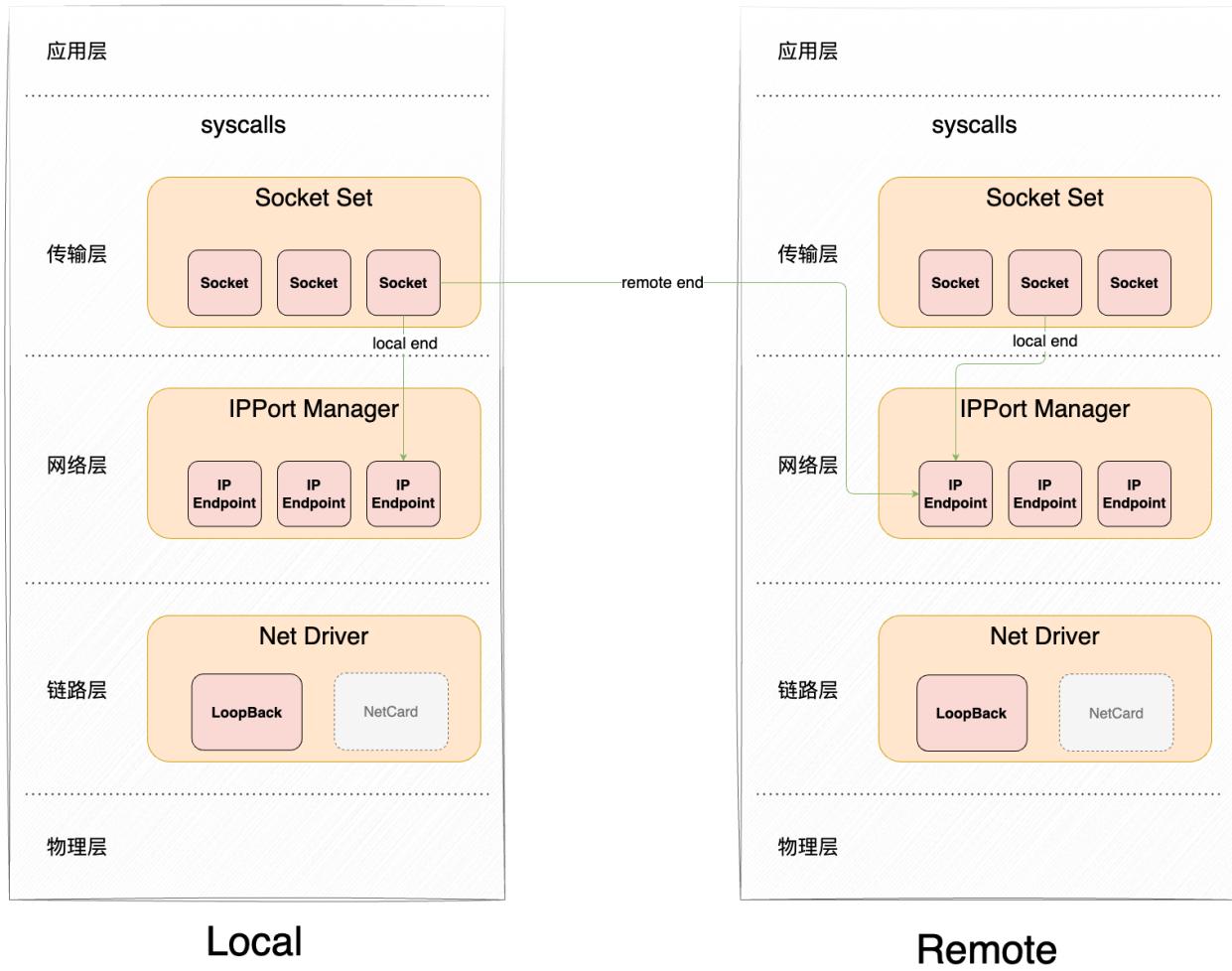


图 29: 网络结构

我们首先设计了 `Socket trait` 用来抽象不同类型套接字的行为 (后文会有这些 `trait` 实现的详细阐述), 当然大部分操作使用了异步让权方法以提高 CPU 利用率:

```
1 /// TCP/UDP or other socket should implement this trait
2 #[async_trait]
3 pub trait Socket: Send + Sync + DowncastSync {
4     /// Read data from the socket.
5     ///
6     /// `buf` is the buffer to store the read data
7     ///
8     /// return:
9     /// - Success: (Returns the length of the data read, the endpoint
10    ///   from which data was read).
11    /// - Failure: Error code
```

```

12     async fn read(&self, buf: &mut [u8]) -> (SysResult<usize>, Option<IpEndpoint>);
13
14     /// Write data to the socket.
15     ///
16     /// `buf` is the data to be written
17     /// `to` is the destination endpoint. If None, the written data will be
18     /// discarded.
19     ///
20     /// return: the length of the data written
21     async fn write(&self, buf: &[u8], to: Option<IpEndpoint>) -> SysResult<usize>;
22
23     /// The bind() function is used to associate a socket with a particular IP
24     /// address and port number on the local machine.
25     ///
26     /// return: whether the operation is successful
27     fn bind(&mut self, local: IpEndpoint, fd: usize) -> SysResult<()>;
28
29     /// `backlog` is the maximum length to which the queue of pending
30     /// connections
31     ///
32     /// return: whether the operation is successful
33     fn listen(&mut self, backlog: usize) -> SysResult<()>;
34
35     /// It is used to establish a connection to a remote server.
36     /// When a socket is connected to a remote server,
37     /// the operating system will establish a network connection with the server
38     /// and allow data to be sent and received between the local socket and the
39     /// remote server.
40     ///
41     /// return: whether the operation is successful
42     async fn connect(&mut self, remote: IpEndpoint) -> SysResult<()>;
43
44     /// It is used to accept a new incoming connection.
45     ///
46     /// only used for TCP
47     async fn accept(&mut self) -> SysResult<(TcpSocket, IpEndpoint)>;
48
49     /// It is used to send data to a connected socket.
50     ///
51     /// return: whether the operation is successful
52     fn shutdown(&mut self, operation: ShutdownType) -> SysResult<()>;
53
54     /// Get the socket's metadata
55     #[allow(unused)]
56     fn meta(&self) -> &SocketMeta;
57
58     /// It is used to get the local endpoint of the socket.
59     fn local_endpoint(&self) -> Option<IpEndpoint>;
60
61     /// It is used to get the remote endpoint of the socket.
62     fn peer_endpoint(&self) -> Option<IpEndpoint>;
63
64     fn setsockopt(&self, _level: usize, _optname: usize, _optval: &[u8]) -> SysResult
65         <()> {
66         Ok(())
67     }

```

在实现网络的系统调用时，只需要调用以上方法即可。在此基础上，根据 Linux 哲学“万物皆文件”，我们对各类 `Socket` 进行了一层包装，给内核提供了 `SocketFile` 结构：

```

1 pub struct TcpSocket {
2     meta: SocketMeta,
3     state: TcpState,
4     handles: Vec<SocketHandle>,
5     local_endpoint: Option<IpEndpoint>,
6 }
7 pub struct UdpSocket {
8     pub handle: SocketHandle,
9     remote_endpoint: Option<IpEndpoint>,
10    meta: SocketMeta,
11 }
12
13 pub enum Sock {
14     Tcp(TcpSocket),
15     Udp(UdpSocket),
16     // Unix(UnixSocket),
17 }
18
19 /// The file for socket
20 pub struct SocketFile {
21     meta: FileMeta,
22     sock: AsyncMutex<Sock>,
23     type_: PosixSocketType,
24 }
```

类似于前文所述的 EXT4 文件系统，我们通过对实际文件包裹了一层 `AsyncMutex`，让其支持在异步环境下的互斥访问。然后为它实现了 `File trait`，包括普通的读写以及支持 I/O 多路复用功能，从而以非常高效的开发效率接入了文件系统，并且也不需要改动任何的常规读写系统调用。

```

1 #[async_trait]
2 impl File for SocketFile {
3     fn meta(&self) -> &FileMeta {
4         &self.meta
5     }
6     async fn base_read(&self, _offset: usize, buf: &mut [u8]) -> SyscallResult {
7         let mut sock = self.socket().await;
8         let res;
9         match &mut *sock {
10             Sock::Tcp(socket) => {
11                 res = socket.read(buf).await.0?;
12             }
13             Sock::Udp(socket) => {
14                 res = socket.read(buf).await.0?;
15             }
16         }
17         Ok(res as isize)
18     }
19     async fn base_write(&self, _offset: usize, buf: &[u8]) -> SyscallResult {
20         let mut sock = self.socket().await;
21         let res;
```

```

22     match &mut *sock {
23         Sock::Tcp(socket) => {
24             res = socket.write(buf, None).await?;
25         }
26         Sock::Udp(socket) => {
27             res = socket.write(buf, None).await?;
28         }
29     }
30     Ok(res as isize)
31 }
32 ...
33 fn poll(&self, req: &PollEvent, waker: Waker) -> PollEvent {
34     let mut sock = block_on(self.socket());
35     let poll_res = match &mut *sock {
36         Sock::Tcp(socket) => socket.poll(req, waker),
37         Sock::Udp(socket) => socket.poll(req, waker),
38     };
39
40     let mut res = PollEvent::empty();
41     if poll_res.contains(PollEvent::POLLIN) {
42         res |= PollEvent::POLLIN;
43     }
44     if poll_res.contains(PollEvent::POLLOUT) {
45         res |= PollEvent::POLLOUT;
46     }
47     if poll_res.contains(PollEvent::POLLHUP) {
48         warn!("[Socket::poll] PollEvent is hangup");
49         res |= PollEvent::POLLHUP;
50     }
51     res
52 }
53 }
```

## 9.2 套接字池 SOCKET\_SET

```

1 pub struct SocketSet {
2     inner: SpinLock<SmoltcpSocketSet<'static>>,
3 }
4 lazy_static::lazy_static! {
5     pub static ref SOCKET_SET: SocketSet = SocketSet::new();
6 }
```

`SOCKET_SET` 是用 `SpinLock` 包裹的 smoltcp 库中的 `SocketSet` 的全局实例。它记录了网络模块中所有的套接字，包括 `Tcp` 和 `Udp`。根据 RAII 的思想，套接字的生命周期就绑定在了该套接字池中，移除时便释放资源。除此之外，它同样维护了所有套接字的状态，在 `LoopBack` 网卡的 `poll` 方法中进行更新。

### 9.2.1 套接字状态更新

借鉴于成熟的操作系统 DragonOS<sup>15</sup>，我们设计了函数 `poll_ifaces` 用来更新套接字的状态：

---

<sup>15</sup><https://github.com/DragonOS-Community/DragonOS>

```

1 lazy_static::lazy_static! {
2     pub static ref NET_DEVICES: RwLock<BTreesMap<usize, Arc<&'static dyn NetworkDev>>> =
3         {
4             let net_devices = RwLock::new(BTreesMap::new());
5             net_devices.write().insert(0, driver::get_net_dev());
6             net_devices
7         };
8
9     pub fn poll_ifaces() {
10         let devices = NET_DEVICES.read();
11         let mut sockets = SOCKET_SET.lock();
12         for (_, iface) in devices.iter() {
13             iface.poll(&mut sockets).ok();
14         }
15     }

```

这个函数会被频繁用于各种操作中，获取最新的套接字状态。

### 9.2.2 套接字生命周期

在任意的套接字被创建时，它会被加入 `SOCKET_SET`：

```

1 let new_socket = Self::new_socket();
2 let new_socket_handle = SOCKET_SET.insert(new_socket);

```

在任意的套接字生命周期结束时，他会从 `SOCKET_SET` 被移除，在移除后通过 `poll_ifaces()` 更新状态，用于 `TcpSocket` 的一端关闭时通知对端关闭，防止对端卡在 `CloseWait` 状态：

```

1 // e.g. UdpSocket
2 fn drop(&mut self) {
3     let mut sockets = SOCKET_SET.lock();
4     let handle = self.handle;
5     let socket = sockets.get_mut::<udp::Socket>(handle);
6     if socket.is_open() {
7         socket.close();
8     }
9     sockets.remove(handle);
10    drop(sockets);
11    poll_ifaces();
12 }

```

### 9.3 端口管理器 PortManager

```

1 pub struct PortItem {
2     pub fd: usize,
3 }
4
5 pub struct PortManager {
6     pub inner: Box<BTreesMap<u16, PortItem>>,
7 }

```

`TCP_PORT_MANAGER` 和 `UDP_PORT_MANAGER` 均为 `PortManager` 的全局实例。`PortManager` 是网络模块中管理端口的核心组件。

### 9.3.1 端口生命周期

同样根据 RAII 的思想，`PortManager` 在套接字执行 `bind` 操作时注册相应的端口，并在生命周期结束时释放对应的端口，由此提供了最大的端口可用性。

```
1 impl PortManager {
2     pub fn bind_port_with_fd(&mut self, port: u16, fd: usize) -> SysResult<u16> {
3         if let Some(_) = &self.inner.get(&port) {
4             error!("[port_manager] Port {port} is already listened (with fd {fd})");
5             Err(Errno::EADDRINUSE)
6         } else {
7             self.inner.insert(port, PortItem::new(fd));
8             Ok(port)
9         }
10    }
11
12    pub fn unbind_port(&mut self, port: u16) {
13        if let Some(_) = self.inner.remove(&port) {
14            debug!("[port_manager] Unbind port {port}");
15        } else {
16            warn!("[port_manager] Port {port} is not listened");
17        }
18    }
19    pub fn resolve_port(&self, endpoint: &IpEndpoint) -> SysResult<u16> {
20        if endpoint.addr.is_unspecified() {
21            let port = if endpoint.port == 0 {
22                self.get_ephemeral_port()?;
23            } else {
24                endpoint.port
25            };
26            Ok(port)
27        } else {
28            assert_ne!(endpoint.port, 0);
29            Ok(endpoint.port)
30        }
31    }
32
33 // e.g. Tcpsocket
34 fn bind(&mut self, local: IpEndpoint, fd: usize) -> SysResult<()> {
35     let mut port_manager = TCP_PORT_MANAGER.lock();
36     let port = port_manager.resolve_port(&local)?;
37     port_manager.bind_port_with_fd(port, fd).unwrap();
38     self.local_endpoint = Some(IpEndpoint::new(local.addr, port));
39     Ok(())
40 }
41 fn drop(&mut self) {
42     if let Some(local) = self.local_endpoint {
43         let mut port_manager = TCP_PORT_MANAGER.lock();
44         port_manager.unbind_port(local.port);
45         drop(port_manager);
46     }
}
```

### 9.3.2 端口复用

在测试集中，出现了多个 `UdpSocket` 绑定到同一个端口的行为，经过测试发现，在远端 `UdpSocket` 读写该端口时，smoltcp 库的行为是仅建立与一个 `UdpSocket` 的连接，其余的套接字得不到任何的响应，这就导致了死锁。

鉴于 smoltcp 库的限制，我们采取了复用文件的逻辑去维护端口复用的逻辑。我们需要保证在所有的时刻，`SOCKET_SET` 中不可能存在多个同种套接字绑定到同样的端口。当一个系统调用尝试操控 `UdpSocket` 绑定到一个已绑定的端口，我们就会复制已绑定端口的 `UdpSocket File` 到当前 `fd` 槽位。

```

1 pub async fn sys_bind(&self, sockfd: usize, addr: usize, addr_len: usize) ->
2     SyscallResult {
3     ...
4     let res = socket.bind(sock_addr, sockfd);
5     match res {
6         Err(Errno::EADDRINUSE) => {
7             // get the old socket file
8             let old_fd = crate::net::get_old_socket_fd(sock_addr.get_endpoint().port);
9
10            // copy the old socket file to the current one, and the current one will be
11            // dropped
12            let mut fd_table = self.task.fd_table();
13            fd_table.copyfrom(old_fd, sockfd)?;
14            drop(fd_table);
15            Ok(0)
16        }
17        r => {
18            r.unwrap();
19            Ok(0)
20        }
21    }
}

```

## 9.4 TcpSocket

TCP (Transmission Control Protocol) 是一种面向连接的、可靠的、基于字节流的传输层协议。在传统的网络套接字实现中，因为缺乏异步无栈协程架构的支持，加之网络传输的特殊性，网络数据的传输延迟以 ms 计，网络模块中存在比别的模块更加严重的忙等现象。于是我们设计了一系列基于无栈协程架构的高效的网络异步操作，在并发场景下，拥有优秀的性能。

首先，对于读写操作，相较于传统的 `loop` 轮询忙等策略而言，我们使用 `yield_now().await` 即时让权，避免了 CPU 时间的浪费，并且利用我们高效的调度器，主动让权大概率可以在下一次被调度时完成读写操作。这样的操作还有一个好处，在高并发情况下因为不涉及 CPU 忙等，CPU 会一直处于高效工作的状态，整体的读写效率会很高。

```

1 // e.g. read
2 async fn read(&self, buf: &mut [u8]) -> (SysResult<usize>, Option<IpEndpoint>) {
3     loop {

```

```

4     poll_ifaces();
5     ...
6     if socket.may_recv() {
7         match socket.recv_slice(buf) {
8             Ok(size) => {
9                 if size > 0 {
10                     ...
11                     return (Ok(size), Some(remote_endpoint.unwrap()));
12                 }
13             }
14             Err(..) => {
15                 return (Err(Errno::ENOTCONN), None);
16             }
17         }
18     } else {
19         return (Err(Errno::ENOTCONN), None);
20     }
21     yield_now().await;
22 }
23 }
```

鉴于 `TcpSocket` 是面向连接的，它的 `connect` 方法是需要被实现的，并且需要维护三次握手的逻辑。我们在这里同样采用了 `yield_now().await`，去屏蔽掉那些忙等时间。并且通过合适的 `poll_ifaces()` 方法，及时的更新套接字状态通知对端连接已经建立，使得连接双方都成为 `Established` 状态。

```

1  async fn connect(&mut self, remote: IpEndpoint) -> SysResult<()> {
2      ...
3      local_socket
4          .connect(iface_inner.context(), remote, temp_port)
5          .map_err(...)?;
6      assert_eq!(local_socket.state(), tcp::State::SynSent);
7
8      loop {
9          poll_ifaces();
10         let mut sockets = SOCKET_SET.lock();
11         let local_socket = sockets.get_mut::<tcp::Socket>(self.handles[0]);
12
13         /*
14          Client                               Server
15
16          CLOSED                                LISTEN
17          |                                     |
18          | ----- send SYN ----->    |
19          SynSent                                |
20          |                                     |
21          |       <---- receive SYN+ACK ----> SYN_RCVD
22          |                                     |
23          | ----- send ACK ----->    |
24          ESTABLISHED                           ESTABLISHED
25        */
26        match local_socket.state() {
27            tcp::State::Closed => {
28                ...
29                yield_now().await;
30            }
31            tcp::State::SynSent => {
```

```

32         ...
33         yield_now().await;
34     }
35     tcp::State::Established => {
36         return Ok(());
37     }
38     _ => {
39         return Err(Errno::ECONNREFUSED);
40     }
41 }
42 }
43 }
```

特殊说明一下，其实 `connect` 时，当前客户端的 `state` 是不可能为 `Closed` 的。此分支只可能说明对端尚未处于监听状态，为防止可能出现的用户程序“客户端先连接，服务端后监听”的情况，也为了增强鲁棒性，我们会在连接失败后，尝试重新连接（也即代码中省略号部分逻辑），当然也使用了让权避免浪费 CPU 时间。

## 9.5 UdpSocket

UDP (User Datagram Protocol)，即用户数据报协议，是一种用于网络通信的传输层协议。它传输速度很快，但可靠性较低。

因为 `UdpSocket` 是面向无连接的，所以我们不关心它的 `accept/listen` 方法。对于读写操作，同样也是类似于 `TcpSocket`，我们采用了 `yield_now().await` 实现了等待让权功能。

```

1 // e.g. read
2 async fn read(&mut self, buf: &mut [u8]) -> (SysResult<usize>, Option<IpEndpoint>) {
3     loop {
4         poll_ifaces();
5         let mut sockets = SOCKET_SET.lock();
6         let socket = sockets.get_mut::<udp::Socket>(self.handle);
7
8         if socket.can_recv() {
9             if let Ok((size, metadata)) = socket.recv_slice(buf) {
10                 ...
11                 return (Ok(size), Some(metadata.endpoint));
12             }
13         }
14         ...
15         yield_now().await;
16     }
17 }
```

对于 `bind` 和 `connect` 操作，其核心在于绑定本地端口和设置远程端点，并且我们对于用户不规范的端点输入，均采取了绑定到随机端口的策略，增加了 `UdpSocket` 的鲁棒性：

```

1 fn bind(&mut self, local: IpEndpoint, fd: usize) -> SysResult<()> {
2     let mut port_manager = UDP_PORT_MANAGER.lock();
3     let port = port_manager.resolve_port(&local)?;
4     let port = port_manager.bind_port_with_fd(port, fd)?;
5     drop(port_manager);
6
7     let mut sockets = SOCKET_SET.lock();
```

```
8     let socket = sockets.get_mut::<udp::Socket>(self.handle);
9
10    if local.addr.is_unspecified() {
11        socket.bind(port)
12    } else {
13        socket.bind(local)
14    }
15    .map_err(...)?;
16
17    Ok(())
18}
19
20 async fn connect(&mut self, remote: IpEndpoint) -> SysResult<()> {
21     self.remote_endpoint = Some(remote);
22     ...
23     if local.port == 0 {
24         let mut port_manager = UDP_PORT_MANAGER.lock();
25         let temp_port = port_manager.get_ephemeral_port()?;
26         port_manager.bind_port(temp_port)?;
27         drop(port_manager);
28         socket.bind(temp_port).map_err(|_| Errno::EINVAL)?;
29     }
30     ...
31     Ok(())
32 }
```

## 9.6 git 测试

目前已经成功完成预发布内容。

```
要添加当前会话，点击左侧的箭头按钮。
1 VisionFive +  
show      Show various types of objects  
status     Show the working tree status  
  
grow, mark and tweak your common history  
backfill   Download missing objects in a partial clone  
branch    List, create, or delete branches  
commit    Record changes to the repository  
merge     Join two or more development histories together  
rebase    Reapply commits on top of another base tip  
reset     Reset current HEAD to the specified state  
switch   Switch branches  
tag       Create, list, delete or verify a tag object signed with GPG  
  
collaborate (see also: git help workflows)  
fetch     Download objects and refs from another repository  
pull      Fetch from and integrate with another repository or a local branch  
push      Update remote refs along with associated objects  
  
'git help -a' and 'git help -g' list available subcommands and some  
concept guides. See 'git help <command>' or 'git help <concept>'  
to read about a specific subcommand or concept.  
See 'git help git' for an overview of the system.  
hint: Using 'master' as the name for the initial branch. This default branch name  
hint: is subject to change. To configure the initial branch name to use in all  
hint: of your new repositories, which will suppress this warning, call:  
hint:  
hint:   git config --global init.defaultBranch <name>  
hint:  
hint: Names commonly chosen instead of 'master' are 'main', 'trunk' and  
hint: 'development'. The just-created branch can be renamed via this command:  
hint:  
hint:   git branch -m <name>  
hint:  
hint: Disable this message with "git config set advice.defaultBranchName false"  
Initialized empty Git repository in /glibc/.git/  
error: cannot run maintenance: No such file or directory  
[master (root-commit) 0bfd0c9] add README.md  
  1 file changed, 1 insertion(+)  
  create mode 100644 README.md  
commit 0bfd0c913af7d5bcd6215bb81988835ea7c3349d (HEAD -> master)  
Author: Your Name <you@example.com>  
Date:   Thu Jan 1 00:01:06 1970 +0000  
  
add README.md  
#### OS COMP TEST GROUP END git-glibc ####  
[kernel] system shutdown (normal exit)
```

图 30: 决赛预发布 git 测例

## 十、物理板卡支持

决赛阶段，我们在进行物理上板的过程中遇到了诸多困难，其中大部分问题都是具体 CPU 型号及其架构相关的硬件特性导致的。在此，我们将上板过程中遇到的架构层相关问题列举出来，以供参考。

### 10.1 LS2K1000 地址未对齐异常

#### ► 编译期地址对齐

通过板卡信息可知，ls2k1000 板卡不支持地址非对齐的访存读写操作，当执行非对齐访存时将会报地址未对齐 (ALE) 异常。对于内核而言，当未打开地址对齐编译选项时，访问栈上地址时会频繁触发 ALE 异常，从而使得性能严重下降。

为了保证内核运行的稳定高效，我们通过关闭部分 LA 下的编译选项来避免异常：

- `-ual` 关闭内存非对齐代码生成，编译器将严格按照对齐地址进行数据维护。
- `-lsx -lasx` 关闭向量扩展指令集优化，2k1000 不支持向量扩展
- `-lvz` 关闭虚拟化支持，2k1000 不支持虚拟化扩展

```
1 [target.loongarch64-unknown-linux-gnu]
2 linker = "loongarch64-linux-gnu-gcc"
3 rustflags = [
4     "-Ctarget-feature=-unal,-lsx,-lasx,-lvz",
5     "-Clink-arg=-Tlib/.ld/linker-la.ld",
6     "-Clink-arg=-nostdlib",
7     "-Clink-arg=-static",
8 ]
```

#### ► 执行期未对齐处理

我们通过编译选项避免了内核态发生的未对齐异常，但用户程序仍然会发生该异常。为了处理用户访存地址未对齐异常，我们在异常处理函数中对于指令进行软件译码并进行单字访存以解决未对齐问题。

由于未对齐指令需要对于程序上下文进行直接操作（读取源寄存器 `rj` 与 `rk`，并修改目标寄存器 `rd`），因此我们修改了初赛时的陷入逻辑，在发生内核陷入时也同样向异常处理程序提供上下文指针，并且设置陷入时首先跳转架构给定的 `<arch>_trap_handler` 中，若检测到未对齐异常则直接处理，然后才进入上层内核进行进一步的判断。

此外该异常处理程序还有多个细节需要注意：

- 与 `qemu` 不同，2k1000 上设置的异常处理入口会自动抹平低 12 位，应当在汇编中规定将其按页对齐，否则会导致异常入口发生错误偏移
- 虽然内核中开启了访存对齐的选项，但仍然无法完全保证栈上信息对齐，选择在栈上保存上下文信息之前，将其存入 `CSR_SAVE1` 中并按 8 字节对齐更为稳妥
- 我们在栈上保存的寄存器上下文超出了原有的 32 寄存器范围，除了通用寄存器外还需要保存 `epc` 等信息，因此将 `sp` 的移动幅度提升到 512 字节以预留充足的空间

```

1 .globl    __kernel_trapvec
2 .align    12
3 __kernel_trapvec:
4 csrwr    $sp, CSR_SAVE1          # save sp since we will align it
5 csrrd    $sp, CSR_SAVE1
6 addi.d   $sp, $sp, -64*8        # reserve enough space for trap frame (currently
7     use 512 bytes)
8 srlid.d  $sp, $sp, 3           # align sp to avoid unaligned exception
9 sllid.d  $sp, $sp, 3
10
11 STORE_X %1
12 STORE_X %2
13 .set      n, 4
14 .rept    28
15 STORE_X %n
16 .set      n, n+1
17 .endr
18
19 csrrd    $t0, CSR_PRMD
20 csrrd    $t1, CSR_ERA
21 csrrd    $t2, CSR_SAVE1        # use SAVE1 to represent sp, we've already saved
22     it
23 st.d     $t0, $sp, 32*8
24 st.d     $t1, $sp, 33*8
25 st.d     $t2, $sp, 3*8         # store sp (x3 in la64)
26
27 move    $a0, $sp              # pass trap frame to handler
28 bl      la_kernel_trap_handler # handle other traps
29
30 ld.d     $t0, $sp, 32*8
31 ld.d     $t1, $sp, 33*8
32 csrwr    $t0, CSR_PRMD       # PRMD
33 csrwr    $t1, CSR_ERA        # ERA
34 LOAD_X %1
35 LOAD_X %2
36 .set      n, 4
37 .rept    28
38 LOAD_X %n
39 .set      n, n+1
40 .endr
41 ld.d     $sp, $sp, 3*8
42 ertn

```

## 10.2 多核启动问题

在进行多核启动时，我们注意到 VisionFive 板卡会固定使用核心 1 而非核心 0 进行启动，这导致我们原有的多核启动维护策略发生了偏差。目前的维护策略是使用充足的 hart 槽位进行 hart 注册初始化，以防止多核核心编号干扰 `per_cpu` 信息的维护，并对于不同板卡使用不同的 `CPU_MASK` 以规范启动时的初始化过程。

```

1 /// CPU count for multi-core
2 pub const CPU_NUM: usize = 4;
3 pub fn wake_other_hart(forbid_hart_id: usize) {

```

```

4     let entry = (_entry_other_hart as usize) & (!KERNEL_ADDR_OFFSET);
5     for i in 0..CPU_NUM {
6         if i != forbid_hart_id && (1 << i) & VALID_PHYS_CPU_MASK != 0 {
7             Arch::hart_start(i, entry);
8         }
9     }
10 }

```

### 10.3 块设备驱动支持

对于 LA 架构，我们对于板卡上搭载的 ahci m2 sata 设备进行了驱动适配，将官方提供的驱动程序接入了 NoAxiom 系统中，并成功启动了块设备。相较于往届 LA 赛道选手实现了真正的块设备访问。

由于良好的硬件抽象设计，我们轻松地实现了驱动的 platform 接口。仅需调用 HAL 层提供的内存抽象层与架构抽象层中的方法，即可完成对于具体硬件功能的实现。同时通过该方法，即使 riscv 架构下的板卡搭载了同类设备，我们理论上仍然能够适配，这实现了设备驱动与底层架构的解耦。

```

1 pub fn ahci_mdelay(ms: u32) {
2     const MS_PER_SEC: usize = 1000;
3     let time = Arch::get_time();
4     let freq = Arch::get_freq();
5     let expire = time + (freq * ms as usize) / MS_PER_SEC;
6     while Arch::get_time() < expire {}
7 }
8 pub fn ahci_sync_dcache() {
9     Arch::sync_dcache();
10 }
11 pub fn ahci_malloc_align(size: u64, align: u32) -> u64 {
12     extern crate alloc;
13     use alloc::boxed::Box;
14     let frame = memory::frame::frame_alloc().unwrap();
15     let frame = Box::leak(Box::new(frame));
16     let addr = frame.kernel_vpn().as_va_usize() as u64;
17     addr
18 }
19 pub fn ahci_phys_to_uncached(pa: u64) -> u64 {
20     pa | IO_ADDR_OFFSET as u64
21 }
22 pub fn ahci_virt_to_phys(va: u64) -> u64 {
23     va & !(KERNEL_ADDR_OFFSET as u64) & !(IO_ADDR_OFFSET as u64)
24 }

```

而对于官方驱动提供的 ahci\_device 结构体，我们将相关代码全部 Rust 风格化重构，并将相关读写接口接入 AhciDevice 的 trait BlockDevice 实现当中，借此向内核暴露统一的设备驱动访问接口。具体代码可见 lib/driver\_ahci。

```

1 impl BlockDevice for LsAhciDevice {
2     fn sync_read(&self, id: usize, buf: &mut [u8]) -> DevResult<usize> {
3         let blknr = id as u64;
4         let blkcnt = (buf.len() / BLOCK_SIZE) as u32;
5         let buffer = buf.as_ptr() as *mut u8;
6         let res = self.device.ahci_sata_read_common(blknr, blkcnt, buffer);
7         match res {

```

```
8     0 => Err(Errno::EIO),
9     _ => Ok(buf.len()),
10    }
11  }
12 fn sync_write(&self, id: usize, buf: &[u8]) -> DevResult<usize> {
13     let blknr = id as u64;
14     let blkcnt = (buf.len() / BLOCK_SIZE) as u32;
15     let buffer = buf.as_ptr() as *mut u8;
16     let res = self.device.ahci_sata_write_common(blknr, blkcnt, buffer);
17     match res {
18         0 => Err(Errno::EIO),
19         _ => Ok(buf.len()),
20     }
21 }
22 }
```

## 十一、参考文献

- **rCore**—*Tutorial*. Available at: <https://rcore-os.cn/rCore-Tutorial-Book-v3/index.html>
- **Pantheon**—*Croutine, Process, Memory*. Available at: [https://gitee.com/LiLiangF/pantheon\\_visionfive](https://gitee.com/LiLiangF/pantheon_visionfive)
- **DragonOS**—*Net structure*. Available at: <https://github.com/DragonOS-Community/DragonOS>
- **Phoenix**—*Croutine, Memory, Signal*. Available at: <https://gitlab.eduxiji.net/educg-group-22026-2376550/T202418123993075-1053>
- **NPUcore-IMPACT**—*LoongArch support for arch*. Available at: <https://gitlab.eduxiji.net/educg-group-22027T202410699992496-1562>
- **Tornado-OS**—*Async driver for file system*. Available at: <https://github.com/HUST-OS/tornado-os>
- **Polyhal**—*Arch layer design for multi-arch*. Available at: <https://github.com/Byte-OS/polyhal>