

基于 RISC-V 架构与Rust语言的16KiB页面操作系统

——操作系统 Project 报告

1 项目简介

- 1.1 项目协作
- 1.2 git 仓库管理
- 1.3 为什么页面大小选择为16KiB，而不是传统甚至是默认的4KiB？
 - 1.3.1 从缓存（按需分页）的角度来看
 - 1.3.2 从虚拟化的角度来看
 - 1.3.3 从分段的角度来看
 - 1.3.4 从硬件接口的角度来看
 - 1.3.5 从内存管理开销(overhead)的角度来看
 - 1.3.6 从TLB的角度来看
- 1.4 为什么选择用Rust语言做Project？
 - 1.4.1 生产力。
 - 1.4.2 可靠性。
 - 1.4.3 高性能。

2 应用程序与基本执行环境

- 2.1 16KiB下的QEMU
- 2.2 Rust 执行环境的实现
 - 2.2.1 实现 致命错误处理器，蓝屏机制
 - 2.2.2 实现log输出
 - 2.2.3 实现条件编译（Makefile+rust feature）
 - 2.2.4 Qemu远程gdb

3 16KiB下内存布局与地址空间的设计

- 3.1 地址转换模式
- 3.2 物理内存分配算法（堆与物理页帧）
 - 3.2.1 best\first\worst
 - 3.2.2 buddy system

4 实现多道程序、分时多任务与进程。

- 4.1 fork() exec() 系统调用, 用户shell
- 4.2 initproc 进程
- 4.3 elf 加载

5 16KiB 相比 4KiB 性能测试

- 5.1 16KiB
- 5.2 4KiB

1 项目简介

1.1 项目协作

姓名	学号	负责内容
叶臻铭	12011404	阅读RISC-V手册，进行QEMU与地址空间设计。用Rust实现并测试物理页分配、进程切换算法。实验测定16KiB方案的优劣。
邓值仁	12012029	掌握Rust语言。深入理解并文档化rCore，深度优化 rCore 设计，使得每个部分都比较稳定可靠、简单清晰、符合rust规范。

1.2 git 仓库管理

本文档的完整最新pdf版[JumbuckNucleus/Project Report.pdf at ch5虚拟内存+进程+ 稳定 · 2catycm/JumbuckNucleus \(github.com\)](https://github.com/2catycm/JumbuckNucleus/blob/master/Project%20Report.pdf)

仓库地址: [2catycm/JumbuckNucleus at ch5虚拟内存+进程+ 稳定 \(github.com\)](https://github.com/2catycm/JumbuckNucleus)

持续集成测试: [feat\(os\): 最新版本 · 2catycm/JumbuckNucleus@896c3a7 \(github.com\)](https://github.com/2catycm/JumbuckNucleus/actions/workflows/feat(os):%20最新版本.yml)

1.3 为什么页面大小选择为16KiB，而不是传统甚至是默认的4KiB?

内存的容量近年来不断发展，容量较大，然而分页方案的页面大小，仍然是上个世纪1960年代就使用的4KiB大小。

¹ 那么，基于16KiB页面的OS有什么优势呢？在做这次Project的之前和过程中，我们需要理解16KiB操作系统的具体优势，这样才能达到Project的效果与目的；否则，就只是做了任务要求而已。

1.3.1 从缓存（按需分页）的角度来看

从按需分页(Demand paging)的角度来看，内存是磁盘的缓存。所谓分页，就是缓存中的分 block，也就是一个传输单位。同一个block或者page的地址(访问需求)，由于访问的内容处于硬件上的同一个传输单位，因此只需第一次访问便将内容缓存，加速大量的后续访问。

那么，是不是block越大，miss rate就越低呢？其实并不尽然。在计算机组成原理的经典教材《计算机组成与设计——硬件/软件接口》中，有这样一幅图。在图中，比较高的曲线的cache容量(按需分页中的内存大小)低，整体失效率比低的曲线要大；横向的变化趋势表示相同内存大小下选择不同的block 大小，失效率的变化。可以看出，如果相对于cache本身的大小，一个block的大小过大，实际上失效率会逐渐上升。²

对于16KiB和4KiB的页面问题来说，这个视角给我们带来了两个启发

- 4KiB页面的 Miss Penalty 较小。即 Page Replacement 完成所需时间较少。
- 16KiB页面虽然Miss Penalty大四倍。但是由于现代计算机的内存容量远远大于4KiB和16KiB，应当适用于图中比较低的曲线，所以16KiB应当显著降低Miss Rate。

1.3.2 从虚拟化的角度来看

上面我们考虑了“物理内存不够，需要磁盘来补充”的早期分页想法之一，然而如今物理内存未必小于虚拟内存。另一个分页的想法来自于虚拟化为不同进程带来的保护。那么，我们应当注意到，不同进程的地址空间不同。

首先考虑进程之间地址空间的隔离

- 假设进程A, B, C, D分别连续访问了4KiB个字节

- 由于四个进程的地址空间不同，他们分别申请了四个不同的16KiB页面而无法互相利用，产生了内碎片。
- 假设进程A连续访问了16KiB个字节，
 - 我们使用16KiB的分页方案，只有一次compulsory miss。
 - 这一次就非常合理地利用了内存空间。

可见，使用多大的页面，和各种进程平均需要访问的内容多少有关系。

另一方面，考虑进程间的通信，进程的内存共享是最简单的方式之一。

- 而内存共享的单位是页。
- 如果页的大小足够，那么通信之间发生的争端(resolution)就会少很多。¹

1.3.3 从分段的角度来看

过大的block会导致内存碎片化的加剧，但是另一方面也减少了内存碎片化。

从外部碎片化的角度来看，

- 分页机制本身作为一种**定长**的分段方法（以前的分段方法是根据程序的code\stack\data\bss这样分的），本身就是为了避免外部碎片化。
 - 定长的页面大小让连续内存分配的算法更加容易设计，因此有更好的性能保证。
- 另一方面，更大的页让内存的可申请尺度变大，
 - 换句话说，碎片的最小大小是一个页，所以不存在比一个页还要小的外部碎片，也就不会出现环节了外部碎片化。

从内部碎片化的角度来看，

- 大的页面导致了内部碎片化。用户可能只申请了一两个字节，就获得了整个物理页面。
- 这种情况叫做系统颠簸(Trashing)。¹

然而，Weisberg 和 Wiseman 认为，**随着内存容量的急剧增大，trashing的问题并不是什么问题**。如果用户需要很大的内存，那么大的页面很合适，如果用户所需内存很小，最多也就浪费一个页面的大小。相比于16GiB的内存，一个进程至多浪费4KiB显得根本不是问题，所以改成16KiB也不会造成什么影响。¹

1.3.4 从硬件接口的角度来看

传统的硬盘设备的一个 Sector 往往就是4KiB，也就是说磁盘一次只能支持 4KiB 大小的传输，这就限制了操作系统的设计也将页面大小设计为4KiB。

然而，近年来，硬盘设备技术发展，16KiB的传输大小并不是什么难事。

1.3.5 从内存管理开销(overhead)的角度来看

- 大的页面，减少了相同大小内存下页面的数量，也就降低了PTE在页表中的数量，进而**减少单个页表的大小**¹也减少了多级页表的页表总数量
- 然而实际上，由于MMU硬件的限制，比如RV64的Sv39方案，
 - MMU硬件只根据PPN，而不是PA来寻找页表。
 - 换句话说，一个页表必须占据一个页的大小。
 - 因此无法有效发挥这一优势。
- 虽然MMU有限制，第二个优点：降低了多级页表的总数量，还是有所体现的。

63	54	53	28	27	19	18	10	9	8	7	6	5	4	3	2	1	0
Reserved			PPN[2]		PPN[1]		PPN[0]		RSW	D	A	G	U	X	W	R	V
10			26		9		9		2	1	1	1	1	1	1	1	1

- 如图为Sv39分页方案，我们称PPN[2:0]对应的页表为二级、一级、零级页表。
- 一个巨页为2MiB, 吉页为1GiB。
- 若连续申请T个字节，至多需要 $1 + \lceil T/2MiB \rceil + \lceil T/1GiB \rceil$ 个页表（二级页表、一级页表、零级页表）³
- 如果改为我们的16KiB页面方案（9,9,9,14）
 - 一个巨页为8MiB, 吉页为4GiB。
 - 那么至多需要 $1 + \lceil T/8MiB \rceil + \lceil T/4GiB \rceil$ 个页表，页表的数量少了大约4倍。

1.3.6 从TLB的角度来看

Weisberg 和 Wiseman 认为，以上因素都有道理，但是之所以16KiB应当是最优选择，最重要的因素是比起4KiB，16KiB可以显著提高TLB的命中率。¹

不妨设内存总大小一定、TLB能存的PTE大小一定、页表能存储的页表项一样多。

假设 TLB能存储256项PTE，Page Size是4KiB。

- 那么TLB能管辖的内存范围是1MiB，如果有连续1MiB的多次访问，TLB在访问过一次后就能对这1MiB范围的所有虚拟地址命中。
- 而如果增大Page Size为16KiB，管辖范围就能提高4倍。
- 对于相同的地址访问来说，TLB 的 miss rate就可以大幅降低。
- 当然，使用巨页也可以增加管辖范围，不过仍然有4倍的关系，而且16KiB优势更为显著。

近年来，硬件设备的速度提升、内存的急剧增大，使得仍为4KiB的页面大小降低了TLB相对于整个内存的覆盖率，因此，提高Page的大小是合理而有必要的。¹

根据 Weisberg 和 Wiseman 的实验结果，16KiB的分页方案比起4KiB降低了1/3的TLB miss数，增加了2/3的内存空间使用，对于如今大的内存硬件而言，是值得的。¹

1.4 为什么选择用Rust语言做Project?

本次 Project 我们没有使用Lab中的C语言代码，而是根据rCore的指导重新构建出一个操作系统。Rust是一门赋予每个人构建可靠且高效软件能力的语言，具有**高性能、可靠性、生产力**这三大优势。⁴本次我们使用Rust开发操作系统Project的过程中，充分运用了这三大优势。

1.4.1 生产力。

1. 可拔插的语言特性。

比起C语言和C++写操作系统没有任何标准库，Rust语言在设计时便考虑了裸机环境、嵌入式设备下没有标准库的问题。使得我们在编写操作系统的时候，仍然可以使用core核心库。

此外，只要在具体平台上编写相应的语言特性，便可以逐步恢复标准库的功能。这一点我们将在后文“Rust执行环境的实现”一节描述。

2. 优质的包管理器。

C/C++的包管理一直以来比较复杂、不能跨平台、不够统一。而Rust作为新兴语言，设计之初就给出了官方的包管理器 `cargo`，统一了调库的接口。

例如，本次Project中我们使用 `riscv` 库和 `log` 库，只需要在项目配置文件 `cargo.toml` 中写：

```
[dependencies]

riscv = { git = "https://github.com/rcore-os/riscv", features = ["inline-asm"] }

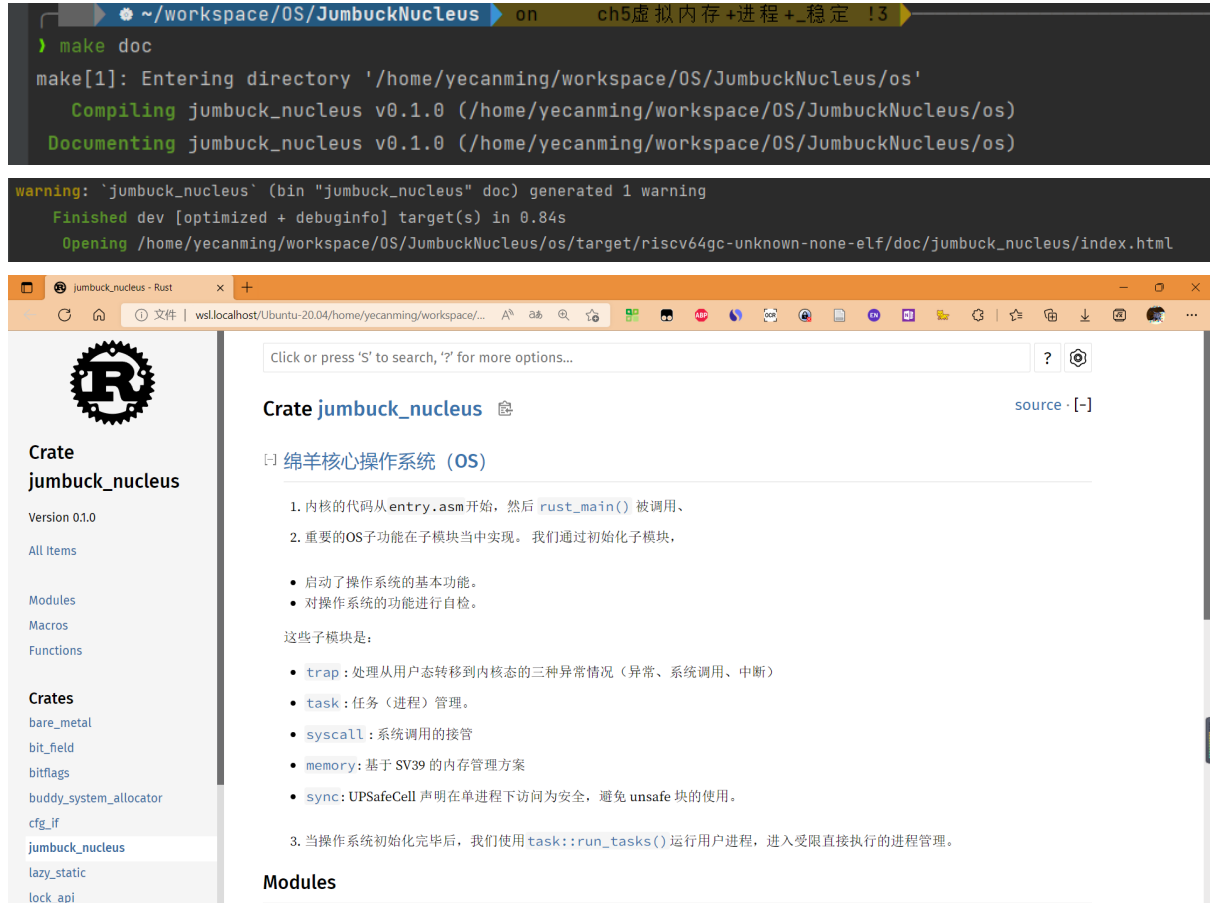
log = "0.4.8"
```

在使用 `cargo build` 构建的时候，便会自动下载对应版本的库(crate)，进行链接。

3. 出色的文档

在Rust中，可以使用`///`的语法为函数、模块等进行文档注释，支持 `markdown` 语法。

接下来，使用 `cargo doc` 可以生成文档，并且在浏览器中查看效果。



可以看到，除了支持markdown语法，rust文档还

- 支持链接别的模块、函数等的文档；
- 形成的文档也是非常美观而统一的，支持简单的搜索和查看源代码的功能
- 支持文档中给出代码使用示例。（图片未体现）

我们查看所调用的库时也可以看它的文档，快速理解用法。

4. 清晰的错误提示

```
warning: function cannot return without recursing
--> src/bin/stack_overflow.rs:7:1
|
7 | fn f(d: usize) {
|   ^^^^^^^^^^^^^^^ cannot return without recursing
8 |     println!("d = {}", d);
9 |     f(d + 1);
|     ----- recursive call site
|
= note: `#[warn(unconditional_recursion)]` on by default
= help: a `loop` may express intention better if this is on purpose
```

这是我们为了测试操作系统的地址空间管理是否可以阻止用户态程序Stack Overflow并且杀死该进程而编写的一个递归调用不结束程序（一直增加函数栈的大小，超出地址空间合法范围）。

Rust编译器针对我们代码给出的一个递归的警告。可以看出，Rust编译器给的信息非常完善，**包括错误的位置、错误的原因、提醒错误原因的笔记和解决错误的帮助。**

而如果我们使用C++来写操作系统，可能使用模板的时候由于编译器历史遗留原因，产生大量不可以理解的奇怪错误。

1.4.2 可靠性。

Rust 丰富的类型系统和所有权模型保证了内存安全和线程安全，让您在编译期就能够消除各种各样的错误。⁴

Rust的管理方法与C++有异曲同工之妙，Rust默认使用C++中称为移动构造的技术，而复制构造则需要实现 `Copy` Trait, 这样移动的语义（称为所有权模型）减少了潜在的复制降低性能的问题。另一方面，Rust和C++都支持RAII的技术，即栈上的资源管理变量（比如Vector,String）拥有指针，在函数控制流比较复杂、甚至包括异常处理的情况下，仍然能够保证资源管理变量能够在合适的时机调用析构函数，自动释放资源。这大量地减少C语言编程中出现的内存泄露、二次释放等问题，而且比Java等GC语言垃圾回收的方案高效很多。

在本次Project中，我们再管理物理页帧的时候使用了RAII这一技术。物理页帧是实现虚拟内存机制的重要部分，实现了物理内存的分页机制后我们才能进一步通过页表管理虚拟内存。操作系统是一种硬件资源管理器。⁵，作为资源管理器，操作系统需要讲物理页帧分配给进程，在进程变为僵尸进程，然后最终被回收的时候，将物理页帧回收。这个过程如果使用C语言编写，可能需要注意不要忘记回收、不要重复回收。如果使用Rust或者C++，我们可以这样写：

```
//mod.rs in frame_allocator. designed by rCore Team and improved by 叶璨铭 to
support continuous allocation.
pub struct FrameBlockTracker {
    pub(super)ppn: PhysPageNum,
    pub(super)count: usize,
}
impl FrameBlockTracker {
    ///创建空的 `FrameTracker`
    pub fn new(ppn: PhysPageNum, count: usize) -> Self {
        // 清理页的内容，全部变为0.
        for i in 0..count{
            PhysPageNum::from(usize::from(ppn)+i).clear();
        }
    }
}
```

```

    }
    Self { ppn, count }
}
}

impl Drop for FrameBlockTracker { //实现析构函数
    fn drop(&mut self) {
        dealloc_frames(self.ppn, self.count);
    }
}

/// 分配一系列连续页面。
/// 接口暴露：本模块外的接口允许调用alloc方法，但是不允许调用dealloc。因为使用
/// [`FrameBlockTracker`]来管理生命周期。
#[allow(unused)]
pub fn alloc_frames(count: usize) -> Option<FrameBlockTracker> {
    FRAME_ALLOCATOR
        .exclusive_access()
        .alloc(count)
        .map(|ppn| FrameBlockTracker::new(ppn, count))
}

fn dealloc_frames(ppn: PhysPageNum, count: usize) {
    FRAME_ALLOCATOR
        .exclusive_access().dealloc(ppn, count);
}

```

首先我们实现了一个 `FRAME_ALLOCATOR`，实现了 `best fit` `worse fit` `buddy system allocator` `first fit` 等连续资源分配算法。但是 `FRAME_ALLOCATOR` 是 `private` 的，我们仅对外暴露 `alloc_frames` 这一接口，外部方法可以通过它获得一个 `FrameBlockTracker` (如果获取失败，则为 `None`)。操作系统其他部分的代码可以通过所有权的转移，延长 `FrameBlockTracker` 的生命周期，而最终析构时，自动调用这里编写的 `drop` 方法，准确而高效地释放相应的内存。

1.4.3 高性能。

Rust 和 C++ 语言一样，提供了“零开销”的抽象。我们用 Rust 写操作系统，不仅不会比 C 语言慢，还可能因为更加清晰的类型系统和所有权机制使得导致变慢的操作在编译时就消除，从而获得更好的内存利用率。

2 应用程序与基本执行环境

2.1 16KiB 下的 QEMU

我们修改了 QEMU 中关于分页的部分。

我们QEMU编译的效果放在了[2catycm/qemu-bin \(github.com\)](https://github.com/2catycm/qemu-bin).

```
) qemu-as-7.0-16Ki
正在设置Qemu为 7.0-16Ki
当前脚本位置 /home/yecanming/workspace/OS/JumbuckNucleus/qemu-bin
设置完成
```

将这个函数集成进makefile中，我们就可以比较方便的切换两种OS。

2.2 Rust 执行环境的实现

2.2.1 实现 致命错误处理器，蓝屏机制

前面提到rust可拔插的语言特性。只要我们实现了 `panic_handler`（致命错误，或者在其他语言中说的不可以处理的异常）属性，那么标准库中各种算法、库报错的时候，都会自动转接到这个函数，我们实现蓝屏信息的打印提示用户，并且可以告知用户错误的位置，最后关机。

```
/// 打印 panic 信息并 [`sbi::shutdown`].
///
/// ### '#[panic_handler]' 属性
/// 声明此函数是 panic 的回调。
#[panic_handler]
fn panic_handler(info: &core::panic::PanicInfo) -> ! {
    console::open_blue_print();
    let repeated_space = "

";

    println!(":\n你的电脑遇到问题，需要关机。");
    println!("我们只为您收集关键信息，然后为您关机。");
    println!("\\t如果您需要了解更多信息，请查看此错误:");
    if let Some(location) = info.location() {
        println!(
            "\\t{:} \\n\\t\\t'{:?}'",
            location.file(),
            location.line(),
            info.message().unwrap()
        );
    }
    } else {
        println!("{}", info.message().unwrap());
    }
    console::close_console_effects();
    sbi::shutdown()
}
```

实现了这个功能之后，我们的OS获得了处理致命错误的能力。

例如，我们试图不断地申请堆内存


```
[jumbuck_nucleus::memory::heap_allocator]
```

```
WARN: 正在执行崩溃测试，接下来操作系统应当崩溃并且关机，而不是死机。
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x80ba0150
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x80b63980
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x80aa6e80
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x8086a380
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x8092d880
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x809f0d80
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x804b4280
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x80577780
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x8063ac80
```

```
[jumbuck_nucleus::memory::heap_allocator]
```

```
INFO: 正在访问物理地址 0x8063ac80
```

```
20%
```

```
30%
```

```
40%
```

```
如果您需要了解更多信息，请查看此错误：
```

```
src/memory/heap_allocator.rs:12
```

```
'绵羊核心堆内存申请异常，可能存在内存泄露！Layout { size_: 4, align_: 4 (1 << 2) }'
```

```
50%
```

```
60%
```

```
70%
```

```
80%
```

```
90%
```

```
100%
```

2.2.2 实现log输出

- 方便debug调试，我们需要一个强大的log输出机制。

2.2.3 实现条件编译（Makefile+rust feature）

- 允许快速选择4KiB还是16KiB

2.2.4 Qemu远程gdb

Qemu支持

3 16KiB下内存布局与地址空间的设计

3.1 地址转换模式

3.2 物理内存分配算法（堆与物理页帧）

3.2.1 best\first\worst

3.2.2 buddy system

4 实现多道程序、分时多任务与进程。

4.1 fork() exec() 系统调用, 用户shell

4.2 initproc 进程

4.3 elf 加载

5 16KiB 相比 4KiB 性能测试

下面我们实际通过实验测试16KiB操作系统是否实际比4KiB的性能更好。

首先编写一个简单的用户态程序。

```
//! 用户程序 bench_pgfault
//! 测试16KiB分页方案是否有效

#![no_std]
#![no_main]

#[macro_use]
extern crate user_lib;
extern crate alloc;
use alloc::vec::Vec;
use core::mem;
use core::f32;
use core::f64;
use user_lib::getpid;
use user_lib::get_time;
#[no_mangle]
pub fn main() -> i32 {
    println!("分页性能测试");
    let start = get_time();
    let mut v: Vec<i32> = Vec::new();
    let l = 16 * 1024 / mem::size_of::<i32>() / 2;
    // 在一个16KiB中反复读写。
    for i in 0..100 {
        for j in 0..1 {
            // v.push((((j as f32)/2.5 + 1.2)*(-1.25)+3.0) / 3.4) as i32);
            v.push((10 * i + j - (i * j)) as i32);
        }
        v.sort();
    }
    println!("测试结束, 用时{}ms", get_time() - start);
}
```

```
0
```

```
}
```

5.1 16KiB

首先运行一次，让缓存加载起来。

```
INFO:  开启时钟中断。
[jumbuck_nucleus]
INFO:  正在加载 App 。
[jumbuck_nucleus]
INFO:  App 加载完成。开始调度运行。
Jumbuck Nucleus OS 绵羊核心操作系统 [Version 0.1.0]
>> bench_pgfault
分页性能测试
测试结束，用时 79ms
user_shell: 进程 2 以退出码 0 结束。
>> █
```

接下来连续运行三次，我们认为这三次运行时间的平均数是16KiB系统运行该程序的期望时间。

```
[jumbuck_nucleus]
INFO:  App 加载完成。开始调度运行。
Jumbuck Nucleus OS 绵羊核心操作系统 [Version 0.1.0]
>> bench_pgfault
分页性能测试
测试结束，用时 58ms
user_shell: 进程 2 以退出码 0 结束。
>> bench_pgfault
分页性能测试
测试结束，用时 55ms
user_shell: 进程 2 以退出码 0 结束。
>> bench_pgfault
分页性能测试
测试结束，用时 53ms
user_shell: 进程 2 以退出码 0 结束。
>> █
```

5.2 4KiB

我们使用

```
exit
fantastic_text
forkexec
forktest
forktest2
forktest_simple
usertests-simple
yield
*****/
Rust user shell
>> bench_pgfault
分页性能测试
测试结束，用时68ms
Shell: Process 2 exited with code 0
>>
```

参考文献

1. P. Weisberg and Y. Wiseman, “Virtual Memory Systems Should use Larger Pages,” Aug. 2015, pp. 1 – 4. doi: [10.14257/asfl.2015.106.01](https://doi.org/10.14257/asfl.2015.106.01). [↩](#) [↩](#) [↩](#) [↩](#) [↩](#) [↩](#) [↩](#)
2. Patterson and Hennessy, *计算机组成与设计——硬件/软件接口(原书第五版) (RISC-V版)*. [↩](#)
3. “rCore-Tutorial-Book-v3 3.6.0-alpha.1 文档.” <http://rCore-os.cn/rCore-Tutorial-Book-v3/index.html> (accessed Jul. 27, 2022). [↩](#)
4. Rust 程序设计语言.” <https://www.rust-lang.org/zh-CN/> (accessed Jul. 29, 2022). [↩](#) [↩](#)
5. R. H. Arpaci-Dusseau and A. C. Arpaci-Dusseau, *Operating systems: Three easy pieces*, 1.00. Arpaci-Dusseau Books, 2018. [↩](#)